

CS331-HW5-Lukang-Sun

October 12, 2021

p1. In my experiments(see Figure1.), 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 $E[\|x_{\text{last step}} - x_\star\|^2]$. As for baselines, I select full gradient descent and CGD with Sparsification($p = 0.2$).

p2. 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 $x_\star = \text{matrix}([-0.54051746] [-0.26890662])$, $x_0 = \text{matrix}([-0.5] [-0.2])$, I sample 10 points in the last step to estimate $E[\|x_{\text{last step}} - x_\star\|^2]$. I set the compression operator on server as identity operator(means there is no compression).

p3. for (a) and (b), without loss of generality, we only consider when $d = 1$.

(a)

$$E[\mathcal{C}_{int}(t)] = (t - \lfloor t \rfloor)(\lfloor t \rfloor + 1) + (1 - t + \lfloor t \rfloor)(\lfloor t \rfloor) = t$$

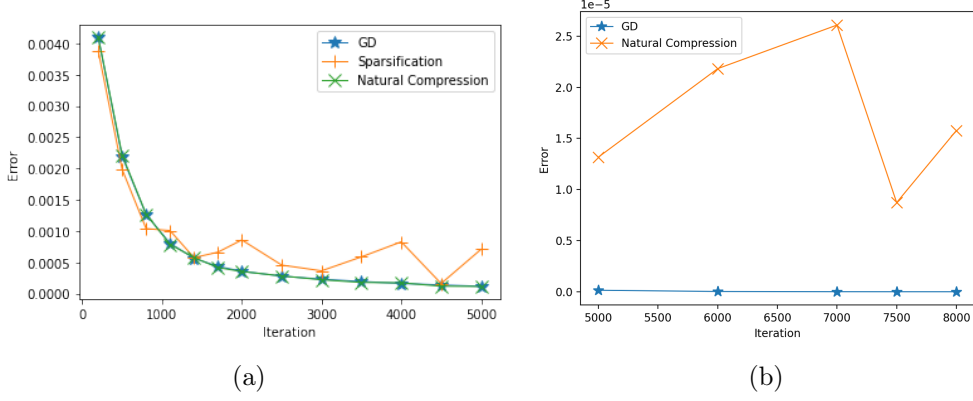


Figure 1: The horizontal axis is iteration, vertical axis is the error. In (a), I set step size as $1/500$, you will see that 1000 steps, full gradient descent and DCGD with natural compression perform quite well, but CGD with sparsification will shake, this can be explained by sparsification operator has bigger ω than natural compression operator. In (b), you will see the difference between DCGD with natural compression and gradient descent, in which, the error of gradient descent keeps converging to 0, while, DCGD with natural compression will shake.

(b) Let's take $t \in (0, 1)$, then $\lfloor t \rfloor = 0$, $\lfloor t \rfloor + 1 = 1$,

$$\mathbb{E} [\|\mathcal{C}_{int}(t)\|^2] = (t - \lfloor t \rfloor)(\lfloor t \rfloor + 1)^2 + (1 - t + \lfloor t \rfloor)(\lfloor t \rfloor)^2 = (t - \lfloor t \rfloor)(\lfloor t \rfloor + 1) + (1 - t + \lfloor t \rfloor)(\lfloor t \rfloor) = t,$$

if $\mathcal{C}_{int} \in \mathbb{B}(\omega)$, this means $t \leq (\omega + 1)t^2$, for any $t \in (0, 1)$, this is impossible, since $1/t \mapsto +\infty$, as $t \mapsto 0$.

(c)

$$\begin{aligned} \mathbb{E} [\|\mathcal{C}_{int}(t) - t\|^2] &= (t - \lfloor t \rfloor)(\lfloor t \rfloor + 1 - t)^2 + (1 - t + \lfloor t \rfloor)(\lfloor t \rfloor - t)^2 \\ &\stackrel{x:=t-\lfloor t \rfloor}{=} x(1-x) \leq \frac{1}{4}, \end{aligned}$$

since $x \in [0, 1)$, the maximum value of $x(1-x)$ is attained at $x = 1/2$, which equals $1/4$. Above we only consider each component, for vector $x \in \mathbb{R}^d$, just take the summation for $i = 1, \dots, d$, we get the upper bound is $\frac{d}{4}$.

p4.

$$\|\text{Top}_k(x) - x\|^2 = \sum_{i=k+1}^d \|x_{\pi_i}\|^2 \quad (1)$$

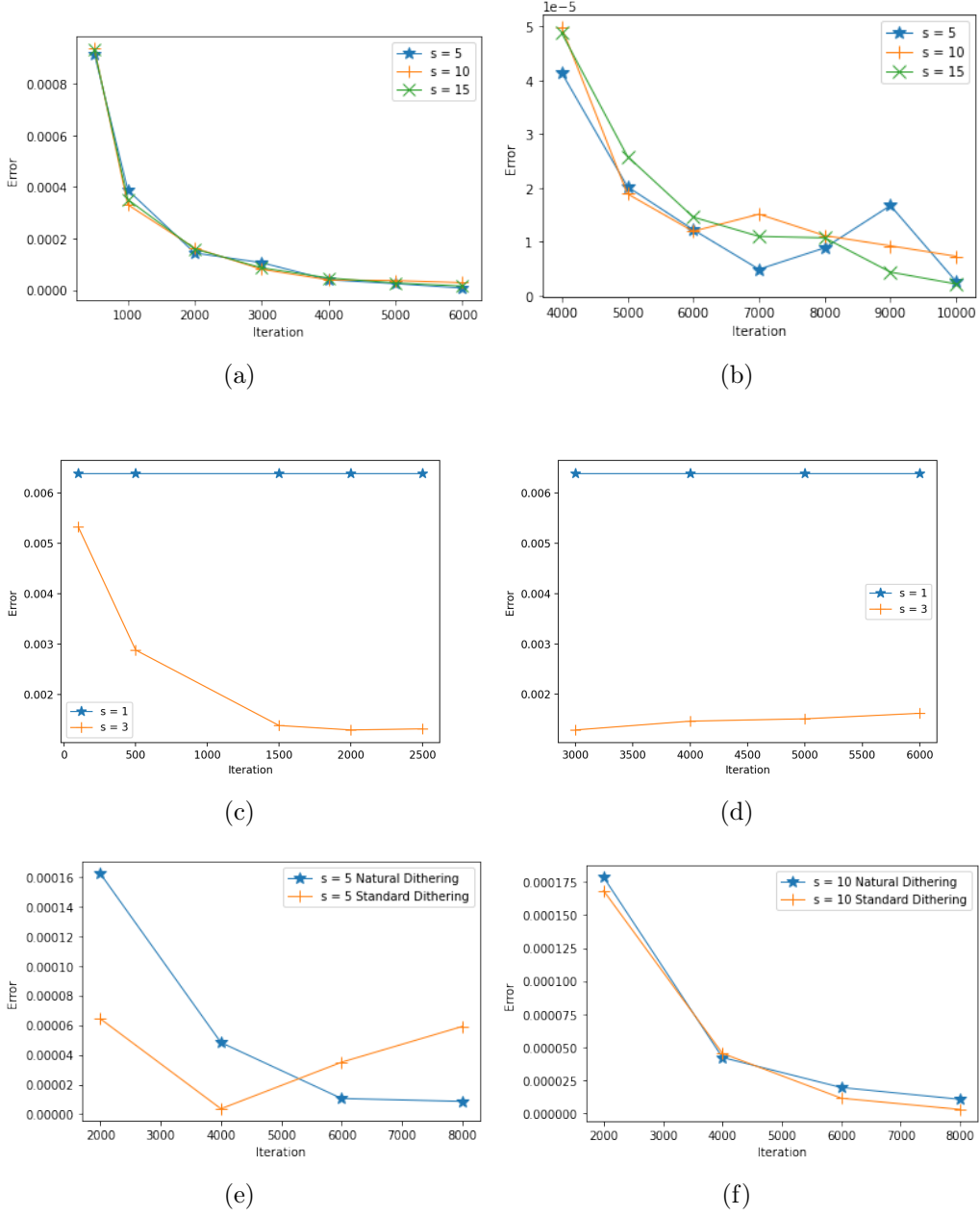


Figure 2: The horizontal axis is iteration, vertical axis is the error. I set step size as $1/500$, in (a) and (b), you will see DCGD with natural dithering with $s = 5, 10, 15$, are quite similar, actually the theory states that $\omega = \frac{1}{8} + \frac{d^{\frac{1}{r}}}{2^{s-1}} \min \left\{ 1, \frac{d^{\frac{1}{r}}}{2^{s-1}} \right\}$, for $s = 5, 10, 15$, $\omega = \frac{1}{8} + \text{term smaller than } 1/16$, they are very close. While in (c) and (d), I set $s = 1, 3$, actually for these two value of s , the corresponding ω are very different and you can also observe that from (c) and (d). In (e) and (f), I show the difference between natural dithering and standard dithering, when s is big (like $s = 10$), the result (f) are similar, while s small ($s = 5$), you can observe the difference.

$$\|x\|^2 = \sum_{i=1}^d \|x_{\pi_i}\|^2 = k\left(\frac{1}{k}\right) \sum_{i=1}^k \|x_{\pi_i}\|^2 + (d-k)\frac{1}{d-k} \sum_{i=k+1}^d \|x_{\pi_i}\|^2 \quad (2)$$

while due to the definition,

$$\frac{1}{k} \sum_{i=1}^k \|x_{\pi_i}\|^2 \geq \frac{1}{d-k} \sum_{i=k+1}^d \|x_{\pi_i}\|^2, \quad (3)$$

so by (2) and (3), we know

$$\|x\|^2 \geq \frac{d}{d-k} \sum_{i=k+1}^d \|x_{\pi_i}\|^2 \stackrel{(1)}{=} \frac{d}{d-k} \|\text{Top}_k(x) - x\|^2, \quad (4)$$

which proves the exercise.