

390015 KU VGSCO: Stochastic Gradient Descent Methods (2022S)

Beryl Ramadhian Aribowo

July 14, 2022

[P4]

E1

Def. 17:

$$D_f(x, y) + D_f(y, x) = \langle \nabla f(x) - \nabla f(y), x - y \rangle = \langle \nabla f(y) - \nabla f(x), y - x \rangle \quad (1)$$

$\forall x, y \in \mathbb{R}^d$:

$$\begin{aligned} \mu \|x - y\|^2 &\leq 2D_f(x, y), \\ \frac{\mu}{2} \|x - y\|^2 &\leq D_f(x, y), \\ \frac{\mu}{2} \|x - y\|^2 &\leq D_f(y, x), \\ D_f(x, y) + \frac{\mu}{2} \|x - y\|^2 &\leq D_f(x, y) + D_f(y, x), \\ D_f(x, y) + \frac{\mu}{2} \|x - y\|^2 &\stackrel{(1)}{\leq} \langle \nabla f(x) - \nabla f(y), x - y \rangle. \end{aligned} \quad (2)$$

E2

$$\begin{aligned} D_f(x, y) + \frac{\mu}{2} \|x - y\|^2 &\leq \langle \nabla f(x) - \nabla f(y), x - y \rangle, \\ \langle \nabla f(x) - \nabla f(y), x - y \rangle &\geq \underbrace{D_f(x, y)}_{\geq \frac{\mu}{2} \|x - y\|^2} + \frac{\mu}{2} \|x - y\|^2, \\ \langle \nabla f(x) - \nabla f(y), x - y \rangle &\geq \frac{\mu}{2} \|x - y\|^2 + \frac{\mu}{2} \|x - y\|^2, \\ \langle \nabla f(x) - \nabla f(y), x - y \rangle &\geq \mu \|x - y\|^2. \end{aligned} \quad (3)$$

[P5]

E4

(i) if $z = x$, then (31) becomes

$$\begin{aligned} D_f(x, x) &= D_f(x, y) + \langle x - y, \nabla f(y) - \nabla f(x) \rangle + D_f(y, x) \\ 0 &= D_f(x, y) + \langle x - y, \nabla f(y) - \nabla f(x) \rangle + D_f(y, x) \\ \langle x - y, \nabla f(x) - \nabla f(y) \rangle &= D_f(x, y) + D_f(y, x). \end{aligned}$$

(ii) if $f(x) = \frac{1}{2} \|x\|^2$, then

$$D_f(x, y) = \frac{1}{2} \|x\|^2 - \frac{1}{2} \|y\|^2 - \langle y, x - y \rangle = \frac{1}{2} \|x - y\|^2,$$

hence (31) becomes

$$\frac{1}{2} \|x - z\|^2 = \frac{1}{2} \|x - y\|^2 + \langle x - y, y - z \rangle + \frac{1}{2} \|y - z\|^2.$$

E5

if $\forall x \in \mathbb{R}^d$

$$\nabla^2 f(x) \succeq 0,$$

then with taylor expansion, for $t \in [0, 1]$

$$f(y) = f(x) + \nabla f(x)^\top (y - x) + \underbrace{\frac{1}{2} (y - x)^\top \nabla^2 f(x + t(y - x)) (y - x)}_{\geq 0}$$

$$f(y) \geq f(x) + \nabla f(x)^\top (y - x),$$

which implies the convexity.

[P6]

E17

(Equation 34):

$$\begin{aligned} \langle a, b \rangle &\leq \frac{\|a\|^2}{2t} + \frac{t\|b\|^2}{2}, \\ \langle a, b \rangle &\leq \frac{\langle a, a \rangle}{2t} + \frac{t\langle b, b \rangle}{2}, \\ 2t\langle a, b \rangle &\leq \langle a, a \rangle + t^2\langle b, b \rangle, \\ 0 &\leq \langle a, a \rangle + \langle tb, tb \rangle - \langle a, tb \rangle - \langle tb, a \rangle, \\ 0 &\leq \|a - tb\|^2. \end{aligned} \tag{4}$$

(Equation 35):

$$\begin{aligned} \|a + b\|^2 &\leq 2\|a\|^2 + 2\|b\|^2, \\ \langle a, a \rangle + \langle b, b \rangle + 2\langle a, b \rangle &\leq 2\langle a, a \rangle + 2\langle b, b \rangle, \\ 0 &\leq \langle a, a \rangle + \langle b, b \rangle - 2\langle a, b \rangle, \\ 0 &\leq \|a - b\|^2. \end{aligned} \tag{5}$$

(Equation 36):

$$\begin{aligned} \frac{1}{2}\|a\|^2 - \|b\|^2 &\leq \|a + b\|^2, \\ \frac{1}{2}\langle a, a \rangle - \langle a, a \rangle &\leq \langle a, a \rangle + \langle b, b \rangle + 2\langle a, b \rangle, \\ \langle a, a \rangle - 2\langle b, b \rangle &\leq 2\langle a, a \rangle + 2\langle b, b \rangle + 4\langle a, b \rangle, \\ 0 &\leq \langle a, a \rangle + \langle 2b, 2b \rangle + \langle a, 2b \rangle + \langle 2b, a \rangle, \\ 0 &\leq \|a + 2b\|^2. \end{aligned} \tag{6}$$

E19

For random vector $X \in \mathbb{R}^d$:

$$\mathbf{Var}[X] := \mathbf{E} [||X - \mathbf{E}[X]||^2]. \quad (7)$$

Markov's inequality:

$$\text{Prob}(X \geq t) \leq \frac{\mathbf{E}[X]}{t}. \quad (8)$$

Proof of Chebyshev's inequality using Markov's inequality:

$$\text{Prob}(|X - \mathbf{E}[X]|^2 \geq t^2) \leq \frac{\mathbf{E}[|X - \mathbf{E}[X]|^2]}{t^2}.$$

Since

$$\text{Prob}(|X - \mathbf{E}[X]|^2 \geq t^2) = \text{Prob}(|X - \mathbf{E}[X]| \geq t), \quad (9)$$

then

$$\text{Prob}(|X - \mathbf{E}[X]| \geq t) \leq \frac{\mathbf{Var}[X]}{t^2}. \quad (10)$$

[P7]

E24

If

$$f = \frac{1}{n} \sum_{i=1}^n f_i,$$

then

$$\begin{aligned}
D_f(x, y) &= \frac{1}{n} \sum_{i=1}^n f_i(x) - \frac{1}{n} \sum_{i=1}^n f_i(y) - \frac{1}{n} \sum_{i=1}^n \langle \nabla f_i(y), x - y \rangle, \\
D_f(x, y) &= \frac{1}{n} \sum_{i=1}^n (f_i(x) - f_i(y) - \langle \nabla f_i(y), x - y \rangle), \\
D_f(x, y) &= \frac{1}{n} \sum_{i=1}^n D_{f_i}(x, y).
\end{aligned}$$

E26

If $\sigma_\star^2 = 0$, then

$$\begin{aligned}
\sigma_\star^2 &= \left(\frac{1}{n^2} \sum_{i=1}^n \frac{||\nabla f_i(x^\star)||^2}{p_i} \right) - ||\nabla f(x^\star)||^2 = 0 \\
&= \left(\frac{1}{n^2} \sum_{i=1}^n \frac{||np_i \nabla f(x^\star)||^2}{p_i} \right) - ||\nabla f(x^\star)||^2 = 0 \\
&= p_i \sum_{i=1}^n (||\nabla f(x^\star)||^2) - ||\nabla f(x^\star)||^2 = 0 \\
&= np_i ||\nabla f(x^\star)||^2 - ||\nabla f(x^\star)||^2 = 0, \\
\sigma_\star^2 = 0 &\implies np_i \nabla f(x^\star) = \nabla f(x^\star).
\end{aligned}$$

E28

For $c > 1$, $i = 1, 2, \dots, n$ and a fixed k ,

$$f_i = \begin{cases} \frac{c}{2} \|x\|^2, & \text{if } i = k \\ \frac{1}{2c} \|x\|^2, & \text{if } i \neq k \end{cases}.$$

In this case, larger c will lead to larger $\max_i L_i$ and larger n will lead to smaller $\frac{1}{n} \sum_i L_i$, hence $\frac{1}{n} \sum_i L_i \ll \max_i L_i$.

[P8]

E33

Let

$$\chi_i = \begin{cases} 1 & i \in S \\ 0 & i \notin S \end{cases}.$$

Since

$$p_i = \frac{1}{n},$$

and

$$|S| = \tau,$$

then

$$\mathbf{E}[\chi_i] = \text{Prob}(i \in S) = \sum_{i=1}^n p_i \chi_i = \frac{1}{n} \sum_{i=1}^n \chi_i = \frac{\tau}{n}.$$

E35

For any vectors, $b_1, \dots, b_n \in \mathbb{R}^d$:

$$\begin{aligned} \left\| \sum_{i=1}^n b_i \right\|^2 - \sum_{i=1}^n \|b_i\|^2 &= \underbrace{\sum_{i=1}^n \langle b_i, b_i \rangle + \sum_{i \neq j} \langle b_i, b_j \rangle}_{\left\| \sum_{i=1}^n b_i \right\|^2} - \sum_{i=1}^n \langle b_i, b_i \rangle, \\ \left\| \sum_{i=1}^n b_i \right\|^2 - \sum_{i=1}^n \|b_i\|^2 &= \sum_{i \neq j} \langle b_i, b_j \rangle. \end{aligned}$$

[P9]

E37

Assumptions of $\mathcal{C} : \mathbb{R}^d \rightarrow \mathbb{R}^d$:

1. $\mathbf{E}[\mathcal{C}(x)] = x, \quad \forall x \in \mathbb{R}^d$
2. $\mathbf{E}[\|\mathcal{C}(x) - x\|^2] \leq \omega \|x\|^2 + \delta, \quad \forall x \in \mathbb{R}^d, \quad \exists \omega, \delta \geq 0$

Proof of convergence for CGD with $n = 1$:
 Since $\mathcal{C} \in \mathbb{B}^d(\omega)$,

$$\mathbf{E} [||g(x)||^2] = \mathbf{E} [||\mathcal{C}(\nabla f(x))||^2] \leq (\omega + 1)||\nabla f(x)||^2. \quad (11)$$

In case of $\nabla f(y) = 0$,

$$\begin{aligned} G(x, y) &:= \mathbf{E} [||g(x) - \nabla f(y)||^2] \\ &= \mathbf{E} [||g(x)||^2] \\ &\stackrel{(11)}{\leq} (\omega + 1)||\nabla f(x) - \nabla f(y)||^2, \\ &\leq 2(\omega + 1)LD_f(x, y). \end{aligned}$$

In case of $\nabla f(y) \neq 0$,

$$\begin{aligned} G(x, y) &:= \mathbf{E} [||g(x) - \nabla f(y)||^2] \\ &= \mathbf{E} [||g(x) - \nabla f(x)||^2] + ||\nabla f(x) - \nabla f(y)||^2 \\ &= \mathbf{E} [||\mathcal{C}(\nabla f(x)) - \nabla f(x)||^2] + ||\nabla f(x) - \nabla f(y)||^2 \\ &\leq \omega||\nabla f(x)||^2 + \delta + ||\nabla f(x) - \nabla f(y)||^2 \\ &= \omega||\nabla f(x) - \nabla f(y) + \nabla f(y)||^2 + ||\nabla f(x) - \nabla f(y)||^2 + \delta \\ &\leq 2\omega||\nabla f(x) - \nabla f(y)||^2 + 2\omega||\nabla f(y)||^2 + ||\nabla f(x) - \nabla f(y)||^2 + \delta \\ &= (2\omega + 1)||\nabla f(x) - \nabla f(y)||^2 + 2\omega||\nabla f(y)||^2 + \delta \\ &\leq 2 \underbrace{(2\omega + 1)L}_{A} D_f(x, y) + \underbrace{2\omega||\nabla f(y)||^2 + \delta}_{C}. \end{aligned}$$

If $0 < \gamma < \frac{1}{A}$, then

$$\mathbf{E} [||x^k - x^*||^2] \leq (1 - \gamma\mu)^k ||x^0 - x^*|| + \frac{2\gamma\omega||\nabla f(x^*)||^2 + \gamma\delta}{\mu}.$$

E39

Lemma 51:

if $\mathcal{C}(x) = x, \forall x$ (no master compression) and $\omega_i = \omega, \forall i$, then

$$G(x, y) \leq 2 \underbrace{\left(L + 2L_{\max} \frac{\omega}{n}\right)}_A D_f(x, y) + \underbrace{2\frac{\omega}{n}\sigma^2(y)}_{C(y)},$$

where

$$\sigma^2(y) := \frac{1}{n} \sum_{i=1}^n ||\nabla f_i(y)||^2.$$

If $\sigma^2(y) = 0$, then

$$G(x, y) \leq 2 \underbrace{\left(L + L_{\max} \frac{\omega}{n}\right)}_A D_f(x, y).$$

Proof:

If $\nabla f(y) \neq 0$, then

$$\begin{aligned} G(x, y) &:= \mathbf{E} [||g(x) - \nabla f(y)||^2] \\ &= \mathbf{E} [||g(x) - \nabla f(x)||^2] + ||\nabla f(x) - \nabla f(y)||^2 \\ &\leq \mathbf{E} [||g(x) - \nabla f(x)||^2] + 2LD_f(x, y), \end{aligned} \quad (12)$$

and

$$g(x) = \mathcal{C}(\hat{g}(x)) = \hat{g}(x) = \frac{1}{n} \sum_{i=1}^n g_i(x). \quad (13)$$

where

$$g_i(x) = \mathcal{C}_i(\nabla f_i(x)).$$

Estimate

$$\begin{aligned} \mathbf{E} [||g(x) - \nabla f(x)||^2] &\stackrel{(13)}{=} \mathbf{E} [||\mathcal{C}(\hat{g}(x)) - \nabla f(x)||^2] \\ &= \mathbf{E} [||\hat{g}(x) - \nabla f(x)||^2] \\ &= \mathbf{E} \left[\left\| \frac{1}{n} \sum_{i=1}^n \underbrace{(g_i(x) - \nabla f_i(x))}_{a_i} \right\|^2 \right] \\ &= \frac{1}{n^2} \mathbf{E} \left[\sum_{i=1}^n ||a_i||^2 + \sum_{i \neq j} \langle a_i, a_j \rangle \right] \\ &= \frac{1}{n^2} \sum_{i=1}^n \mathbf{E} [||a_i||^2] + \sum_{i \neq j} \mathbf{E} [\langle a_i, a_j \rangle] \\ &= \frac{1}{n^2} \sum_{i=1}^n \mathbf{E} [||a_i||^2] + \sum_{i \neq j} \underbrace{\langle \mathbf{E}[a_i], \mathbf{E}[a_j] \rangle}_0 \\ &\leq \frac{1}{n^2} \sum_{i=1}^n \omega_i ||\nabla f_i(x)||^2 \\ &= \frac{\omega}{n^2} \sum_{i=1}^n ||\nabla f_i(x)||^2. \end{aligned}$$

Next, bound

$$\begin{aligned} ||\nabla f_i(x)||^2 &= ||\nabla f_i(x) - \nabla f_i(y) + \nabla f_i(y)||^2 \\ &\leq 2||\nabla f_i(x) - \nabla f_i(y)||^2 + 2||\nabla f_i(y)||^2 \\ &\leq 4L_i D_{f_i}(x, y) + 2||\nabla f_i(y)||^2. \end{aligned}$$

Combine everything:

$$\begin{aligned} G(x, y) &\leq \mathbf{E} [||g(x) - \nabla f(x)||^2] + 2LD_f(x, y) \\ &\leq \frac{\omega}{n^2} \sum_{i=1}^n ||\nabla f_i(x)||^2 + 2LD_f(x, y) \\ &\leq \frac{\omega}{n^2} \sum_{i=1}^n (4L_i D_{f_i}(x, y) + 2||\nabla f_i(y)||^2) + 2LD_f(x, y) \\ &= 2\frac{\omega}{n} \left(2 \sum_{i=1}^n \frac{1}{n} L_i D_{f_i}(x, y) + \frac{1}{n} \sum_{i=1}^n ||\nabla f_i(y)||^2 \right) + 2LD_f(x, y) \\ &\leq 2\frac{\omega}{n} (2L_{\max} D_f(x, y) + \sigma^2(y)) + 2LD_f(x, y) \\ &= 2(L + 2L_{\max})D_f(x, y) + 2\frac{\omega}{n} \sigma^2(y). \end{aligned} \quad (14)$$

Else, if $\nabla f(y) = 0$, then

$$\begin{aligned}
G(x, y) &= \mathbf{E} [||\hat{g}(x)||^2] \\
&= \mathbf{E} [||\hat{g}(x) - \mathbf{E} [\hat{g}(x)]||^2] + ||\mathbf{E} [\hat{g}(x)]||^2 \\
&= \mathbf{E} [||\hat{g}(x) - \nabla f(x)||^2] + ||\nabla f(x)||^2 \\
&\leq \left(\frac{\omega}{n^2} \sum_{i=1}^n ||\nabla f_i(x)||^2 \right) + ||\nabla f(x) - \nabla f(y) + \nabla f(y)||^2 \\
&\leq \left(\frac{\omega}{n^2} \sum_{i=1}^n ||\nabla f_i(x)||^2 \right) + 2||\nabla f(x) - \nabla f(y)||^2 + 2||\underbrace{\nabla f(y)}_0||^2 \quad (15) \\
&\leq \left(\frac{\omega}{n^2} \sum_{i=1}^n ||\nabla f_i(x)||^2 \right) + 2LD_f(x, y) \\
&\leq \dots \text{ same as (14), from the third line} \\
&= 2(L + 2L_{\max})D_f(x, y) + 2\frac{\omega}{n}\sigma^2(y).
\end{aligned}$$

[P10]

E41

Let

$$p_i = \text{Prob}(i \in S),$$

where

$$S \subseteq \{1, 2, \dots, d\},$$

then

$$\mathbf{E}[|S|] = \mathbf{E} \left[\sum_{i=1}^d |S_i| \right] = \sum_{i=1}^d \mathbf{E}[|S_i|] = \sum_{i=1}^d 1p_i + 0(1 - p_i) = \sum_{i=1}^d p_i.$$

E42

If $\mathbf{E}[\mathbf{C}^\top \mathbf{C}]$ is finite, then $\forall x \neq 0$:

$$\begin{aligned}
x^\top \mathbf{E}[\mathbf{C}^\top \mathbf{C}]x &\geq 0 \\
\mathbf{E}[x^\top \mathbf{C}^\top \mathbf{C}x] &\geq 0 \\
x^\top \mathbf{C}^\top \mathbf{C}x &\geq 0 \\
(\mathbf{C}x)^\top (\mathbf{C}x) &\geq 0.
\end{aligned}$$

[P11]

E47

Define base case:

$$C_{1,2} := C_1 \circ C_2 \in \mathbb{B}^d(\underbrace{((\omega_1 + 1)(\omega_2 + 1) - 1)}_{\omega_{1,2}}),$$

$$C_{1,3} := C_1 \circ C_2 \circ C_3 = C_{1,2} \circ C_3,$$

$$\begin{aligned}\omega_{1,3} &:= (\omega_{1,2} + 1)(\omega_3 + 1) - 1 \\ &= ((\omega_1 + 1)(\omega_2 + 1) - 1 + 1)(\omega_3 + 1) - 1 \\ &= (\omega_1 + 1)(\omega_2 + 1)(\omega_3 + 1) - 1,\end{aligned}$$

$$C_{1,n} := C_1 \circ C_2 \circ \dots \circ C_n = C_{1,n-1} \circ C_n,$$

$$\omega_{1,n} := (\omega_1 + 1)(\omega_2 + 1) \dots (\omega_n + 1) - 1 = (\omega_{1,n-1} + 1)(\omega_n + 1) - 1.$$

By induction, the base case is clear. Next, if $n = k$, assume

$$C_{1,k} := C_1 \circ C_2 \circ \dots \circ C_k = C_{1,k-1} \circ C_k,$$

$$\omega_{1,k} := (\omega_1 + 1)(\omega_2 + 1) \dots (\omega_k + 1) - 1 = (\omega_{1,k-1} + 1)(\omega_k + 1) - 1$$

is true. Then for $n = k + 1$:

$$C_{1,k+1} := C_1 \circ C_2 \circ \dots \circ C_k \circ C_{k+1} = C_{1,k-1} \circ C_k \circ C_{k+1},$$

$$\begin{aligned}\omega_{1,k+1} &:= (\omega_{1,k} + 1)(\omega_{k+1} + 1) - 1 = (\omega_1 + 1)(\omega_2 + 1) \dots (\omega_k + 1)(\omega_{k+1} + 1) - 1 \\ &:= ((\omega_{1,k-1} + 1)(\omega_k + 1) - 1 + 1)(\omega_{k+1} + 1) - 1 = (\omega_1 + 1)(\omega_2 + 1) \dots (\omega_k + 1)(\omega_{k+1} + 1) - 1 \\ &:= ((\omega_1 + 1)(\omega_2 + 1) \dots (\omega_k + 1))(\omega_{k+1} + 1) - 1 = (\omega_1 + 1)(\omega_2 + 1) \dots (\omega_k + 1)(\omega_{k+1} + 1) - 1, \\ \omega_{1,k+1} &= (\omega_1 + 1)(\omega_2 + 1) \dots (\omega_k + 1)(\omega_{k+1} + 1) - 1.\end{aligned}$$

E48

For $a, b \in \mathbb{R}^d$, define

$$\min\{a_i, b_i\} = \begin{cases} a_i, & \text{if } a_i < b_i \\ b_i, & \text{if } a_i > b_i \end{cases}.$$

Thus

$$\sum_i \min\{a_i, b_i\} = \begin{cases} \sum_i a_i, & \text{if } a_i < b_i, \forall i \\ \sum_i b_i, & \text{if } a_i > b_i, \forall i \end{cases}.$$

In the case of inequality

$$\min\{a_i, b_i\} < a_i \text{ or } b_i,$$

thus

$$\sum_i \min\{a_i, b_i\} < \sum_i a_i \text{ or } \sum_i b_i.$$

[P12]

E55

The DCGD-SHIFT has the same exact steps in the algorithm as DCGD ($n \geq 1$ case), the difference is the gradient estimator:

$$g_h(x) := \frac{1}{n} \sum_{i=1}^n g_{h_i}(x) = \frac{1}{n} \sum_{i=1}^n h_i + \mathcal{C}_i(\nabla f_i(x) - h_i). \quad (16)$$

which means the gradients on the workers are shifted and then compressed. In order to prove the convergence theorem, first decompose

$$\mathbf{E} [\|g_h(x^k) - \nabla f(x^*)\|^2] = \mathbf{E} [\|g_h(x^k) - \nabla f(x^k)\|^2] + \|\nabla f(x^k) - \nabla f(x^*)\|^2.$$

Then, bound

$$\begin{aligned} \mathbf{E} [\|g_h(x^k) - \nabla f(x^k)\|^2] &= \mathbf{E} \left[\left\| \frac{1}{n} \sum_{i=1}^n \underbrace{\mathcal{C}_i(\nabla f_i(x^k) - h_i) + h_i - \nabla f_i(x^k)}_{b_i^k} \right\|^2 \right] \\ &= \frac{1}{n^2} \mathbf{E} \left[\sum_i \|b_i^k\|^2 \sum_{i \neq j} \langle b_i^k, b_j^k \rangle \right] \\ &= \frac{1}{n^2} \sum_{i=1}^n \mathbf{E} [\|b_i^k\|^2] + \frac{1}{n^2} \sum_{i \neq j} \underbrace{\langle \mathbf{E}[b_i^k], \mathbf{E}[b_j^k] \rangle}_0 \\ &= \frac{1}{n^2} \sum_{i=1}^n \mathbf{E} [\|\mathcal{C}_i(\nabla f_i(x^k) - h_i) + h_i - \nabla f_i(x^k)\|^2] \\ &\leq \frac{1}{n^2} \sum_{i=1}^n \omega_i \|\nabla f_i(x^k) - h_i\|^2 \\ &= \frac{1}{n^2} \sum_{i=1}^n \omega_i \|\nabla f_i(x^k) - \nabla f_i(x^*) - (h_i - \nabla f_i(x^*))\|^2 \\ &\leq \frac{2}{n^2} \sum_{i=1}^n \omega_i \|\nabla f_i(x^k) - \nabla f_i(x^*)\|^2 + \omega_i \|h_i - \nabla f_i(x^*)\|^2 \\ &\leq \frac{2}{n^2} \sum_{i=1}^n 2\omega_i L_i D_{f_i}(x^k, x^*) + \frac{2}{n^2} \sum_{i=1}^n \omega_i \|h_i - \nabla f_i(x^*)\|^2 \\ &\leq \frac{4}{n} \max(L_i \omega_i) \frac{1}{n} \sum_{i=1}^n D_{f_i}(x^k, x^*) + \frac{2}{n^2} \sum_{i=1}^n \omega_i \|h_i - \nabla f_i(x^*)\|^2 \\ &\leq \frac{4}{n} \max(L_i \omega_i) D_{f_i}(x^k, x^*) + \frac{2}{n^2} \sum_{i=1}^n \omega_i \|h_i - \nabla f_i(x^*)\|^2. \end{aligned}$$

Thus

$$\mathbf{E} [\|g_h(x^k) - \nabla f(x^*)\|^2] \leq \underbrace{2 \left(L + \frac{2}{n} \max(\omega_i L_i) \right) D_f(x^k, x^*)}_A + \underbrace{\frac{2}{n^2} \sum_{i=1}^n \omega_i \|h_i - \nabla f_i(x^*)\|^2}_C.$$

[P13]

E56

Thm 94: whenever $B = 0$ and $M = 0$, then $\frac{B+M\tilde{B}}{M} = 0$, proof:

First, the stepsize γ satisfies

$$0 < \gamma < \frac{1}{\mu}. \quad (17)$$

Then the iterates $\{x^k, \sigma^k\}$ satisfy

$$\mathbf{E}[d^k] \leq (1 - \gamma\mu)^k d^0 + \frac{C\gamma}{\mu}. \quad (18)$$

where

$$d^k := \|x^k - x^*\|^2. \quad (19)$$

From Lemma 95, it is clear

$$\mathbf{E}[d^{k+1}] \leq (1 - \gamma\mu)\mathbf{E}[d^k] + C\gamma^2.$$

By recurrence, we obtain

$$\mathbf{E}[d^k] \leq (1 - \gamma\mu)^k d^0 + \frac{C\gamma}{\mu}.$$

[P14]

E57

In case of arbitrary p , the gradient estimator of L-SVRG is

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

Hence, the unbiasedness:

$$\begin{aligned} \mathbf{E}[g^k | x^k, y^k] &= \mathbf{E}[g(x^k) - g(y^k) + \nabla f(y^k) | x^k, y^k] \\ &= \mathbf{E}[g(x^k) | x^k, y^k] - \mathbf{E}[g(y^k) | x^k, y^k] + \mathbf{E}[\nabla f(y^k) | x^k, y^k] \\ &= \nabla f(x^k) - \nabla f(y^k) + \nabla f(y^k) \\ &= \nabla f(x^k). \end{aligned}$$

E58

If

$$g(x) = \nabla f(x) + \xi,$$

then

$$\begin{aligned} g^k &= (\nabla f(x^k) + \xi) - (\nabla f(y^k) + \xi) + \nabla f(y^k) \\ &= \nabla f(x^k), \end{aligned}$$

which is exactly GD's gradient estimator, where in this case p does not have any role, since the gradient estimator does not depend on y^k anymore. The convergence rate in this case, with stepsize $\gamma = \frac{1}{6A''}$ is

$$\mathbf{E}[d^k] \leq \left(1 - \frac{\mu}{6A''}\right)^k d^0,$$

where

$$d^k := \|x^k - x^*\|^2.$$

Thus

$$k \geq \frac{6A''}{\mu} \log \frac{1}{\epsilon},$$

which is equal to GD's rate of $\mathcal{O}\left(\frac{L}{\mu} \log \frac{1}{\epsilon}\right)$.

[P15]

E62

The algorithm with (200) as the update rule is equivalent to CGD if we set

$$\begin{aligned} x^k &:= h_i^k, \\ g^k &:= \mathcal{C}_i^k(\nabla f_i(x^k) - h_i^k). \end{aligned}$$

Corollary 49 says that, if $0 < \gamma \leq \frac{1}{(\omega+1)L}$ and $\nabla f(x^*) = 0$, then

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

In case of one step iteration, then

$$\mathbf{E} \left[\|x^{k+1} - x^*\|^2 \right] \leq (1 - \gamma\mu) \|x^k - x^*\|^2. \quad (20)$$

Since the optimization problem is in the form of

$$\max_{h_i} \phi_i^k(h_i) := -\frac{1}{2} \|h_i - \nabla f_i(x^k)\|^2, \quad (21)$$

then the solution is

$$\nabla \phi_i^k(h_i) = \nabla f_i(x^k) - h_i = 0 \implies \nabla f_i(x^k) = h_i. \quad (22)$$

which corresponds to $\nabla f(x^*) = 0$ in CGD's case. Also, it can be noticed that ϕ_i^k is 1-smooth and 1-strongly convex, i.e., $\mu = L = 1$. Thus, with $0 < \alpha \leq \frac{1}{w_i+1}$, (200) is equivalent to (20).

E63

The update rule

$$h^{k+1} = h^k - \alpha \mathcal{C}(h^k - \nabla f(x^k)) \quad (23)$$

can be interpreted as a *descent* rule instead of *ascent*, and minimize rather than maximize. The solution for the minimization of ϕ_i^k is (22), the only difference is the compressed shifted gradient is

$$\tilde{g}^k := h^k - \nabla f(x^k) = -(\nabla f(x^k) - h^k),$$

where this does not change the convergence properties since $\|\tilde{g}^k\|^2$ is considered for the convergence bounds. Finally, we have

$$h_i^{k+1} = h_i^k - \alpha \mathcal{C}_i(\tilde{g}_i^k)$$

which is the interpretation of descent. Hence (200) still holds.

If DIANA uses (23) as the update rule then the convergence properties of it does not change by similar argument to when (21) uses (23) as update rule. The bound on the AC inequality does not change, since

$$\|\nabla f(x) - h\| = \|h - \nabla f(x)\|,$$

and the bound of the σ^k assumption also does not change.

[P16]

E68

(216) can be transformed into Lagrange multiplier form, which is

$$\mathcal{L}(J, y) := \frac{1}{2} \|\mathbf{J} - \mathbf{J}^k\|_F^2 + y_l^\top (\mathbf{J} e_l - \nabla f_l(x^k)).$$

The solution can be obtained by

$$\begin{aligned} \frac{\partial}{\partial \mathbf{J}_{ij}} \mathcal{L} &= 0 \\ \mathbf{J}_{ij} - \mathbf{J}_{ij}^k + y_{li} e_{lj} &= 0 \\ \mathbf{J}_{ij} &= \mathbf{J}_{ij}^k + y_{li} e_{lj}, \end{aligned} \tag{24}$$

where

$$e_{lj} = \begin{cases} 1 & \text{if } l = j \\ 0 & \text{otherwise} \end{cases}.$$

If the last line of (24) is inserted back to the equality constraint, then

$$\begin{aligned} \text{if } l \neq j, & \text{ then } \mathbf{J}_{ij}^k e_i = \nabla f_i(x^k), \\ \text{else if } l = j, & \text{ then } \mathbf{J}_{ij}^k e_i + y_{li} e_i = \nabla f_i(x^k). \end{aligned}$$

E70

If $\mathbf{S} = e_i$, then (231) becomes

$$\mathbf{J}^{k+1} = \mathbf{J}^k + (\mathbf{J}(x^k) - \mathbf{J}^k) e_i e_i^\top,$$

$(\mathbf{J}(x^k) - \mathbf{J}^k) e_i e_i^\top$ vanishes if $i \neq j$, hence

$$\mathbf{J}_{:,j}^{k+1} = \mathbf{J}_{:,j}^k,$$

else, if $i = j$ then

$$\mathbf{J}_{:,j}^{k+1} = \mathbf{J}(x^k) e_i e_i^\top = \nabla f_i(x^k).$$

[P17]

E71

Since x^\star is an optimal solution of (232), by the definitions in (232), (233), and (234)

$$\begin{aligned} f(x^\star) &= \hat{f}(z) \\ \frac{1}{n} \sum_i f_i(x^\star) &= \frac{1}{n} \sum_i f_i(z_i), \end{aligned}$$

hence

$$z_1 = z_2 = \dots = z_n = x^\star,$$

in other words

$$z^\star := (x^\star, \dots, x^\star) \in \mathcal{L}.$$

is an optimal solution of (233).

Conversely, if $z^* := (z_1^*, \dots, z_n^*) \in \mathbb{R}^{nd}$ is an optimal solution of (233), then

$$\hat{f}(z^*) = \frac{1}{n} \sum_i f_i(z_i^*) = \frac{1}{n} \sum_i f_i(x) = f(x),$$

which means

$$\begin{aligned} x &= z_1^* = z_2^* = \dots = z_n^*, \\ z^* &\in \mathcal{L}, \end{aligned}$$

hence $x^* := z_1^*$ is an optimal solution.

E73

DIANA with control update step

$$\hat{h}^{k+1} \begin{cases} \nabla \hat{f}(z^k) & \text{with probability } p \\ \hat{h}^k & \text{with probability } 1 - p \end{cases}, \quad (25)$$

introduces stochasticity to \hat{h} . Since

$$\begin{aligned} \nabla \hat{f}(z^k) &= \frac{1}{n} (\nabla f_1(z_1^k), \dots, \nabla f_n(z_n^k)) \\ &= \frac{1}{n} (\nabla f_1(x^k), \dots, \nabla f_n(x^k)), \end{aligned} \quad (26)$$

and by multiplying both sides with n , (25) becomes

$$n\hat{h}^{k+1} \begin{cases} (\nabla f_1(x^k), \dots, \nabla f_n(x^k)) & \text{with probability } p \\ n\hat{h}^k & \text{with probability } 1 - p \end{cases},$$

with $b_j^k := n\hat{h}_j^k$ and by using the last line of (26), (25) can be rewritten as

$$b_j^{k+1} \begin{cases} \nabla f_j(x^k) & \text{with probability } p \\ b_j^k & \text{with probability } 1 - p \end{cases}$$

for $j \in [n]$. Meanwhile the gradient estimator used for the induced SGD is the same as (255). Hence the induced SGD method is as the following:

Algorithm 1 DIANA - stocH

- 1: **Parameters:** learning rate $\gamma > 0$; initial iterate $x^0 \in \mathbb{R}^d$; initial control vectors $b_1^0, \dots, b_n^0 \in \mathbb{R}^d$.
- 2: **for** $k = 0, 1, 2, \dots$ **do**
- 3: Sample $i \in [n]$ uniformly at random
- 4: Compute gradient estimate: $g^k = \nabla f_i(x^k) - b_i^k + \frac{1}{n} \sum_{j=1}^n b_j^k$
- 5: Update iterate: $x^{k+1} = x^k - \gamma g^k$
- 6: Update control vectors, $\forall j \in [n]$:

$$b_j^{k+1} \begin{cases} \nabla f_j(x^k) & \text{with probability } p \\ b_j^k & \text{with probability } 1 - p \end{cases}$$

7: **end for**
