

CS331-HW10-Lukang-Sun

November 12, 2021

p1. (see Figure 1.) $a = [\text{matrix}([0.1]), \text{matrix}([0.424466]), \text{matrix}([0.77981303]), \text{matrix}([0.20033184]), \text{matrix}([0.51116473]), \text{matrix}([0.2604399]), \text{matrix}([0.97100656]), \text{matrix}([0.21263449]), \text{matrix}([0.26417151]), \text{matrix}([0.15995097])]$

$b = [\text{matrix}([0.4231786]), \text{matrix}([0.524466]), \text{matrix}([0.17981303]), \text{matrix}([0.50033184]), \text{matrix}([0.71116473]), \text{matrix}([0.0604399]), \text{matrix}([0.37100656]), \text{matrix}([0.91263449]), \text{matrix}([0.66417151]), \text{matrix}([0.65995097])]$, $f = \frac{1}{10} \sum_{i=1}^{10} f_i(x, y)$, $f_i(x, y) = \sin(x + a[i]) + \cos(y + b[i])$, for the SGD method, I use SGD-US and L-SVRG-US. Initial point is $\text{init} = \text{matrix}([-0.5], [-0.2])$.

(1) L-SVRG-US will converge to the optimal point, while SGD-US will only converge to a neighborhood of the optimal point, in graph (a), we can see the results verifies prediction.

(2) in L-SVRG-US, p is bigger the trajectories will be more smooth, while p is smaller, the trajectories will fluctuate, in graph (b), we can see the results verifies prediction.

p2.

Proof. The proof is almost the same as the proof of theorem 121. Since f is L -smooth, we have

$$\begin{aligned} f(x^{k+1}) - f^{\inf} &\leq f(x^k) - f^{\inf} + \langle \nabla f(x^k), x^{k+1} - x^k \rangle + \frac{L}{2} \|x^{k+1} - x^k\|^2 \\ &= f(x^k) - f^{\inf} - \gamma \langle \nabla f(x^k), g^k \rangle + \frac{L\gamma^2}{2} \|g^k\|^2 \end{aligned}$$

By applying expectation to both sides and subsequently using unbiasedness of g^k and the assumed bound on the second moment of the stochastic gradient,

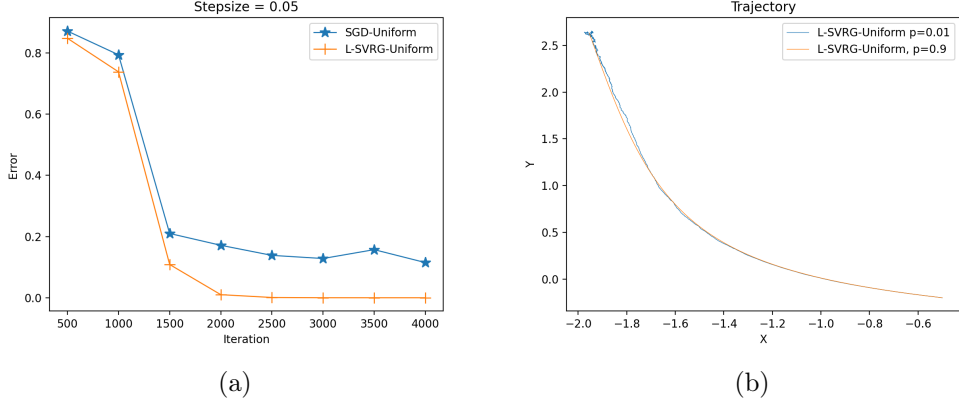


Figure 1: (a) shows $E[\|\nabla f\|^2]$ changes in terms of iteration number K when set step size $\gamma = 0.05$, (b) shows the trajectories of L-SVRG-US with different p .

we get

$$\begin{aligned}
E[f(x^{k+1}) - f^{\inf}] &\leq E[f(x^k) - f^{\inf}] - \gamma E[\|\nabla f(x^k)\|^2] + \frac{L\gamma^2}{2} E[\|g^k\|^2] \\
&\leq E\left[f(x^k) - f^{\inf} - \gamma \|\nabla f(x^k)\|^2 + \frac{L\gamma^2}{2} \left[2A(f(x^k) - f^{\inf}) + B_1\sigma^k + B_2\|\nabla f(x^k)\|^2 + C\right]\right] \\
&= E\left[\left(1 + LA\gamma^2\right)(f(x^k) - f^{\inf}) + \frac{LB_1\gamma^2}{2}\sigma^k - \left(\gamma - \frac{LB_2\gamma^2}{2}\right)\|\nabla f(x^k)\|_2^2 + \frac{LC\gamma^2}{2}\right]
\end{aligned}$$

Choose any $M > 0$ and define

$$\Delta^{k+1} \stackrel{\text{def}}{=} f(x^{k+1}) - f^{\inf} + M\gamma^2\sigma^{k+1}$$

$$\begin{aligned}
E[\Delta^{k+1}] &\leq E\left[\left(1 + LA\gamma^2 + 2M\tilde{A}\gamma^2\right)(f(x^k) - f^{\inf}) + \left(\frac{LB_1}{2} + M\tilde{B}_1\right)\gamma^2\sigma^k\right] \\
&\quad - E\left[\left(\gamma - \frac{LB_2\gamma^2}{2} - M\tilde{B}_2\gamma^2\right)\|\nabla f(x^k)\|^2 + \frac{LC\gamma^2}{2} + M\tilde{C}\gamma^2\right] \\
&= E\left[a\left[f(x^k) - f^{\inf} + \frac{\frac{LB_1}{2} + M\tilde{B}_1}{a}\gamma^2\sigma^k\right] - b\|\nabla f(x^k)\|^2 + c\right]
\end{aligned}$$

where

$$\begin{aligned} a &\stackrel{\text{def}}{=} 1 + LA\gamma^2 + 2M\tilde{A}\gamma^2 \\ b &\stackrel{\text{def}}{=} \gamma - \frac{LB_2\gamma^2}{2} - M\tilde{B}_2\gamma^2 \\ c &\stackrel{\text{def}}{=} \frac{LC\gamma^2}{2} + M\tilde{C}\gamma^2 \end{aligned}$$

In order to turn (1) into a recursion which has Δ^k on the right hand side, we need to make sure that

$$\frac{\frac{LB_1}{2} + M\tilde{B}_1}{a} \leq M$$

Fortunately, it is easy to see (prove this!) that we can make sure this holds by an appropriate choice of M . In particular, the last inequality holds if we choose

$$M \stackrel{\text{def}}{=} \frac{LB_1}{2(1 - \tilde{B}_1)} = \frac{L\theta}{2}$$

With this choice of M , we can continue from (314) and obtain the recursion

$$\mathbb{E} [\Delta^{k+1}] \leq \mathbb{E} [a\Delta^k - b \|\nabla f(x^k)\|^2] + c \quad (1)$$

By applying expectation to both sides of this, and using the tower property of expectation, we get the recursion

$$\begin{aligned} \mathbb{E} [\Delta^{k+1}] &= \mathbb{E} [\mathbb{E} [\Delta^{k+1} \mid x^k, \xi^k]] \\ &\leq a\mathbb{E} [\Delta^k] - b\mathbb{E} [\|\nabla f(x^k)\|^2] + c \end{aligned} \quad (2)$$

We now apply Lemma 120 to recursion (2) with $X_k = \mathbb{E} [\Delta^k]$ and $Y_k = b\mathbb{E} [\|\nabla f(x^k)\|^2]$. If we set $x = x^k$ with probability p_k (where p_k is as in Lemma 120), which means that $Y = Y_k$ with probability p_k , we conclude that

$$\begin{aligned} b\mathbb{E} [\|\nabla f(x)\|^2] &= \mathbb{E}[Y] \\ &\leq \frac{a^K}{S_K} \Delta^0 + c \\ &\leq \frac{a^K}{K} \Delta^0 + c \end{aligned} \quad (3)$$

where the last inequality follows since $a \geq 1$, which implies that $S_K \geq K$.

We now evaluate the expressions for b and c in (3). First,

$$\begin{aligned}
b &= \gamma - \frac{LB_2\gamma^2}{2} - M\tilde{B}_2\gamma^2 \\
&= \gamma - \frac{\gamma}{2} \left(LB_2\gamma + L\theta\tilde{B}_2\gamma \right) \\
&\geq \frac{\gamma}{2}
\end{aligned} \tag{4}$$

where the last inequality holds by setting

$$\gamma \leq \frac{1}{L(B_2 + \theta\tilde{B}_2)}$$

Moreover,

$$c = \frac{LC}{2}\gamma^2 + M\tilde{C}\gamma^2 = \frac{L}{2}(C + \theta\tilde{C})\gamma^2 \tag{5}$$

By plugging the bound (4) on b and expression (5) for c into (3), we obtain the result. \square

p3.

Proof. By lecture, we have

$$E[\|\nabla f(x)\|^2] \leq \frac{2(f(x^0) - f^{\inf})}{K} \times \max\left\{ \underbrace{\sqrt{\frac{4}{3} \frac{1-p}{p} \alpha(c+1)}}_{M_1}, \underbrace{L \left(B_2 + \frac{B_1\tilde{B}_2}{1-\tilde{B}_1} \right)}_{M_2} \right\}$$

$$M_2 \leq L \left(1 + \frac{3}{c+1} + \frac{4(1-p)}{pc} \right)$$

$$\alpha \stackrel{\text{def}}{=} \frac{(n-\tau)L_{\text{avg}}^2}{(n-1)\tau},$$

if we let $p = 0.5, c = 1$, then $M_1 = \sqrt{\frac{8}{3}}n^{-1/3}L_{\text{avg}}, M_2 \leq 6.5L$, so $\max\{M_1, M_2\} \leq 6.5L$, when n is large. So we need $\mathcal{O}\left(\frac{1}{\epsilon^2}\right)$ steps to make the error less than ϵ^2 , each step's computation is $n^{2/3}$, so, the total complexity is $\mathcal{O}(n^{2/3}/\epsilon^2)$. \square

p4.

Theorem. Let Assumption 8 (L -smoothness and $f \geq f^{\inf}$), PL-condition ($\|\nabla f(x)\|_2^2 \geq 2\mu(f(x) - f(x^*))$) and Assumption 11 (σ^k assumption for non-convex functions) be satisfied. Choose constant stepsize γ satisfying

$$0 < \gamma \leq \min \left\{ \frac{1}{\frac{LB_2}{2} + M\tilde{B}_2}, \frac{\mu}{LA + 2MA + \mu LB_2 + 2\mu M\tilde{B}_2} \right\}$$

Then

$$\mathbb{E} [\Delta^k] \leq (1 - \mu\gamma)^k \Delta^0 + \frac{c}{\mu\gamma}$$

where $c \stackrel{\text{def}}{=} \frac{LC\gamma^2}{2} + M\tilde{C}\gamma^2$, $M = \frac{LB_1}{2(a-2\mu b-\tilde{B}_1)}$, $a \stackrel{\text{def}}{=} 1 + LA\gamma^2 + 2M\tilde{A}\gamma^2$, $b \stackrel{\text{def}}{=} \gamma - \frac{LB_2\gamma^2}{2} - M\tilde{B}_2\gamma^2$, $\Delta^k \stackrel{\text{def}}{=} f(x^k) - f^{\inf} + M\gamma^2\sigma^k$,

Proof. Since f is L -smooth, we have

$$\begin{aligned} f(x^{k+1}) - f^{\inf} &\leq f(x^k) - f^{\inf} + \langle \nabla f(x^k), x^{k+1} - x^k \rangle + \frac{L}{2} \|x^{k+1} - x^k\|^2 \\ &\stackrel{(295)}{=} f(x^k) - f^{\inf} - \gamma \langle \nabla f(x^k), g^k \rangle + \frac{L\gamma^2}{2} \|g^k\|^2. \end{aligned} \tag{6}$$

By applying expectation to both sides of (6) and subsequently using unbiasedness of g^k and the assumed bound (297) (from lecture) on the second moment of the stochastic gradient, we get

$$\begin{aligned} \mathbb{E} [f(x^{k+1}) - f^{\inf} \mid x^k, \xi^k] &\stackrel{(6)}{\leq} f(x^k) - f^{\inf} - \gamma \|\nabla f(x^k)\|^2 + \frac{L\gamma^2}{2} \mathbb{E} [\|g^k\|^2 \mid x^k, \xi^k] \\ &\stackrel{(297)}{\leq} f(x^k) - f^{\inf} - \gamma \|\nabla f(x^k)\|^2 \\ &\quad + \frac{L\gamma^2}{2} [2A(f(x^k) - f^{\inf}) + B_1\sigma^k + B_2\|\nabla f(x^k)\|^2 + C] \\ &= (1 + LA\gamma^2)(f(x^k) - f^{\inf}) + \frac{LB_1\gamma^2}{2}\sigma^k \\ &\quad - \left(\gamma - \frac{LB_2\gamma^2}{2} \right) \|\nabla f(x^k)\|^2 + \frac{LC\gamma^2}{2} \end{aligned} \tag{7}$$

Choose any $M > 0$ and define

$$\Delta^{k+1} \stackrel{\text{def}}{=} f(x^{k+1}) - f^{\inf} + M\gamma^2\sigma^{k+1} \tag{8}$$

by combining inequality (7) with assumption (298), we get

$$\begin{aligned}
E [\Delta^{k+1} \mid x^k, \xi^k] &\stackrel{(306)+(298)}{\leq} \underbrace{\left(1 + LA\gamma^2 + 2M\tilde{A}\gamma^2\right)}_{=a} (f(x^k) - f^{\inf}) + \left(\frac{LB_1}{2} + M\tilde{B}_1\right) \gamma^2 \sigma^k \\
&\quad - \underbrace{\left(\gamma - \frac{LB_2\gamma^2}{2} - M\tilde{B}_2\gamma^2\right)}_{=b} \|\nabla f(x^k)\|^2 + \underbrace{\frac{LC\gamma^2}{2} + M\tilde{C}\gamma^2}_{=c} \\
&\stackrel{\text{PL-condition}}{\leq} (a - 2\mu b) \left[f(x^k) - f^{\inf} + \frac{\frac{LB_1}{2} + M\tilde{B}_1}{a - 2\mu b} \gamma^2 \sigma^k \right] + c
\end{aligned} \tag{9}$$

where

$$\begin{aligned}
a &\stackrel{\text{def}}{=} 1 + LA\gamma^2 + 2M\tilde{A}\gamma^2 \\
b &\stackrel{\text{def}}{=} \gamma - \frac{LB_2\gamma^2}{2} - M\tilde{B}_2\gamma^2 \\
c &\stackrel{\text{def}}{=} \frac{LC\gamma^2}{2} + M\tilde{C}\gamma^2,
\end{aligned}$$

we also require $b \geq 0$, that is $\gamma \leq \frac{1}{\frac{LB_2}{2} + M\tilde{B}_2}$. In order to turn (314) into a recursion which has Δ^k on the right hand side, we need to make sure that

$$\frac{\frac{LB_1}{2} + M\tilde{B}_1}{a - 2\mu b} \leq M$$

We can choose $M = \frac{LB_1}{2(a - 2\mu b - \tilde{B}_1)}$. With this choice of M , we can continue from (10) and obtain the recursion

$$E [\Delta^{k+1} \mid x^k, \xi^k] \leq (a - 2\mu b) \Delta^k + c \tag{10}$$

By applying expectation to both sides of this, and using the tower property of expectation, we get the recursion

$$\begin{aligned}
E [\Delta^{k+1}] &= E [E [\Delta^{k+1} \mid x^k, \xi^k]] \\
&\leq (a - 2\mu b) E [\Delta^k] + c
\end{aligned}$$

Finally,

$$E [\Delta^k] \leq (a - 2\mu b)^k \Delta^0 + c \sum_{i=0}^{k-1} (a - 2\mu b)^i \tag{11}$$

by choosing $\gamma \leq \frac{\mu}{LA + 2MA + \mu LB_2 + 2\mu M\tilde{B}_2}$, we will have $a - 2\mu b \leq 1 - \mu\gamma$ so from (11), we have

$$E [\Delta^k] \leq (1 - \mu\gamma)^k \Delta^0 + \frac{c}{\gamma\mu} \tag{12}$$

□