

CS331-HW8-Lukang-Sun

October 29, 2021

p1.

Proof. let $u := \text{prox}_{\gamma R}(x), v := \text{prox}_{\gamma R}(y)$, then by definition we have, $0 \in \partial R(u) + \frac{1}{\gamma}(u - x), 0 \in \partial R(v) + \frac{1}{\gamma}(v - y)$, so $x - u - (y - v) \in \gamma(\partial R(u) - \partial R(v))$, by strong convexity of R , we have

$$\langle x - u - (y - v), u - v \rangle \geq \gamma \lambda \|u - v\|_2^2, \quad (1)$$

so

$$\begin{aligned} \|x - y\|^2 &= \|x - u - (y - v) + (u - v)\|^2 \\ &\geq \|u - v\|^2 + 2\langle x - u - (y - v), u - v \rangle \\ &\geq (1 + 2\gamma\lambda)\|u - v\|^2. \end{aligned} \quad (2)$$

This proves the problem. \square

p2.

Theorem. Assume that f is μ -convex and R is λ -convex, g^k is unbiased (Assumption 1) and that the AC assumption (Assumption 2) is satisfied. Choose a stepsize satisfying

$$0 < \gamma \leq \frac{1}{A}$$

Then the iterates $\{x^k\}_{k \geq 0}$ of SGD (Algorithm 3) satisfy

$$E \left[\|x^k - x^*\|^2 \right] \leq \left(\frac{1 - \mu\gamma}{1 + 2\lambda\gamma} \right)^k \|x^0 - x^*\|^2 + \frac{C\gamma}{\mu(1 + 2\lambda\gamma)}$$

Proof. let $r^k \stackrel{\text{def}}{=} x^k - x^*$. Performing the exact same steps as in the case of

GD, but replacing the gradient $\nabla f(x^k)$ by the stochastic gradient g^k , we get

$$\begin{aligned}
\|r^{k+1}\|^2 &= \|\text{prox}_{\gamma R}(x^k - \gamma g^k) - x^*\|^2 \\
&= \|\text{prox}_{\gamma R}(x^k - \gamma g^k) - \text{prox}_{\gamma R}(x^* - \gamma \nabla f(x^*))\|^2 \\
&\stackrel{(p1.)}{\leq} \frac{1}{1+2\gamma\lambda} \|(x^k - \gamma g^k) - (x^* - \gamma \nabla f(x^*))\|^2 \\
&= \frac{1}{1+2\gamma\lambda} (\|x^k - x^* - \gamma(g^k - \nabla f(x^*))\|^2) \\
&= \frac{1}{1+2\gamma\lambda} (\|r^k\|^2 - 2\gamma \langle r^k, g^k - \nabla f(x^*) \rangle + \gamma^2 \|g^k - \nabla f(x^*)\|^2),
\end{aligned}$$

so

$$\begin{aligned}
\mathbb{E} [\|r^{k+1}\|^2 \mid x^k] &\leq \frac{1}{1+2\gamma\lambda} (\|r^k\|^2 - 2\gamma \mathbb{E} [\langle r^k, g^k - \nabla f(x^*) \rangle \mid x^k] + \gamma^2 \mathbb{E} [\|g^k - \nabla f(x^*)\|^2 \mid x^k]) \\
&= \frac{1}{1+2\gamma\lambda} (\|r^k\|^2 - 2\gamma \langle r^k, \mathbb{E} [g^k \mid x^k] - \nabla f(x^*) \rangle + \gamma^2 \mathbb{E} [\|g^k - \nabla f(x^*)\|^2 \mid x^k]) \\
&= \frac{1}{1+2\gamma\lambda} (\|r^k\|^2 - 2\gamma \langle r^k, \nabla f(x^k) - \nabla f(x^*) \rangle + \gamma^2 \mathbb{E} [\|g^k - \nabla f(x^*)\|^2 \mid x^k]),
\end{aligned}$$

use

$$\langle r^k, \nabla f(x^k) - \nabla f(x^*) \rangle \geq D_f(x^k, x^*) + \frac{\mu}{2} \|x^k - x^*\|^2,$$

lead to

$$\mathbb{E} [\|r^{k+1}\|^2 \mid x^k] \leq \frac{1}{1+2\gamma\lambda} \left((1 - \gamma\mu) \|r^k\|^2 - 2\gamma D_f(x^k, x^*) + \gamma^2 \mathbb{E} [\|g^k - \nabla f(x^*)\|^2 \mid x^k] \right),$$

use Assumption (2), we get

$$\mathbb{E} [\|r^{k+1}\|^2 \mid x^k] \leq \frac{1}{1+2\lambda\gamma} (1 - \gamma\mu) \|r^k\|^2 + \frac{1}{1+2\lambda\gamma} \gamma^2 C,$$

so we have

$$\mathbb{E} [(1 + 2\lambda\gamma)^k \|r^k\|^2] \leq (1 - \mu\gamma)^k \|r^0\|^2 + \frac{\gamma C (1 + 2\lambda\gamma)^{k-1}}{\mu},$$

divide by $(1 + 2\lambda\gamma)^k$ from both sides, we have

$$\mathbb{E} [\|r^k\|] \leq \left(\frac{1 - \mu\gamma}{1 + 2\lambda\gamma} \right)^k \|r^0\|^2 + \frac{C\gamma}{\mu(1 + 2\lambda\gamma)}$$

□

p3. (i)

$$E[g^k | \mathcal{F}_k(y)] = (1 - q)\nabla f(y^k) + q\left(\frac{1}{q}\nabla f(x^k) + \left(1 - \frac{1}{q}\right)\nabla f(y^k)\right) = \nabla f(x^k)$$

(ii)

Lemma. Assume f is μ -convex and L -smooth, R is convex, then we have

$$E\left[\|g^k - \nabla f(x^*)\|^2 | x^k, \xi^k\right] \leq 2AD_f(x^k, x^*) + B\sigma^k,$$

$$E[\sigma^{k+1} | x^k, \xi^k] \leq 2\tilde{A}D_f(x^k, x^*) + \tilde{B}\sigma^k,$$

where $\sigma^k = E[\|\nabla f(y^k) - \nabla f(x^*)\|^2]$, $A = \frac{2L}{q}$, $B = 1 - q + \frac{2(1-q)^2}{q}$, $\tilde{A} = pL$, $\tilde{B} = (1 - p)$.

Theorem. Assume that f is μ -convex. For any $M > \frac{B}{1-\tilde{B}}$ choose a stepsize γ satisfying

$$0 < \gamma \leq \min\left\{\frac{1}{\mu}, \frac{1}{A + M\tilde{A}}\right\}$$

Then the iterates $\{x^k, \sigma^k\}_{k \geq 0}$ of SGD-CTRL satisfy

$$E[V^k] \leq \max\left\{(1 - \gamma\mu)^k, \left(\frac{B + M\tilde{B}}{M}\right)^k\right\} V^0$$

where the Lyapunov function V^k is defined by

$$V^k \stackrel{def}{=} \|x^k - x^*\|^2 + M\gamma^2\sigma^k$$

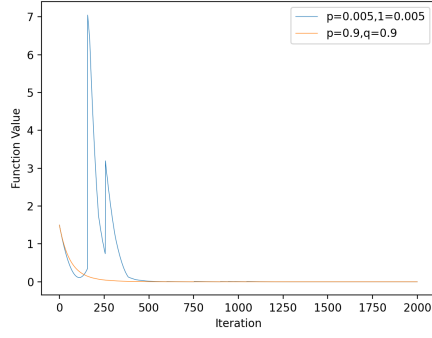
Proof. We only need to prove the lemma, the theorem is a corollary of the lemma.

$$\begin{aligned} E[\|g^k - \nabla f(x^*)\|^2 | k] &= (1 - q)\|\nabla f(x^k) - \nabla f(x^*)\|^2 \\ &\quad + q\|1/q(\nabla f(x^k) - \nabla f(x^*)) + (1 - 1/q)(\nabla f(y^k) - \nabla f(x^*))\|^2 \\ &\leq (1 - q + \frac{2(1-q)^2}{q})\|\nabla f(y^k) - \nabla f(x^*)\|^2 + 2(\frac{2L}{q})D_f(x^k, x^*) \end{aligned}$$

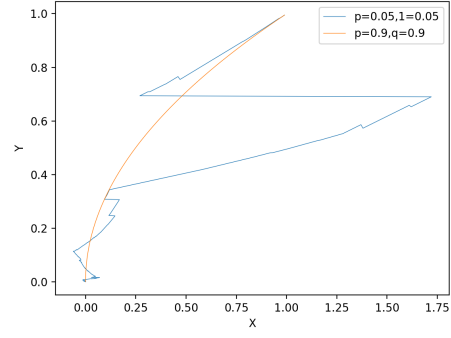
$$\begin{aligned} E[\|\nabla f(y^{k+1}) - \nabla f(x^*)\|^2 | k] &= (1 - p)\|\nabla f(y^k) - \nabla f(x^*)\|^2 + p\|\nabla f(x^k) - \nabla f(x^*)\|^2 \\ &\leq (1 - p)\sigma^k + 2(pL)D_f(x^k, x^*). \end{aligned}$$

□

(iii) when $p \approx 1, q \approx 1$, this method perform well, when $p \approx 0, q \approx 0$, this method performs poorly. This method will reduce the computation of the gradient, so when it is very expensive to compute the gradient (like the dimension is very large), then this method might be a good substitute of gradient descent method. In my experiment, $f = x^2 + 0.5y^2$, step size $\gamma = 0.005$, I get the results in terms of different p, q .



(a)



(b)

Figure 1: (a) shows function value change in terms of iteration, (b) shows the trajectories of this method with different p, q . you can see that when p, q are close to 1, this method converges quite smoothly, when p, q are close to 0, this method converges with fluctuations, this matches the theory.