

# HW3-sol-Optimization in machine learning

Tian QIU

August 28, 2024

## 1 SGD for Smooth and Strongly Convex Functions

### 1.1 a

By  $l$ -smoothness:

$$\begin{aligned}\mathbb{E}[f(x_{s+1}) - f(x_s)] &\leq \mathbb{E}\left[\langle \nabla f(x_s), x_{s+1} - x_s \rangle + \frac{l}{2} \|x_{s+1} - x_s\|^2\right] \\ &= \mathbb{E}[\langle \nabla f(x_s), -\eta g(x_s) \rangle] + \frac{\eta^2 l}{2} \mathbb{E}[\|g(x_s)\|^2] \\ &= -\eta \nabla f(x_s)^T \mathbb{E}[g(x_s)] + \frac{\eta^2 l}{2} \mathbb{E}[\|g(x_s)\|^2].\end{aligned}$$

By the property and assumption of variance:

$$\mathbb{E}[\|g(x_s)\|^2] - \|\mathbb{E}[\nabla f(x_s)]\|^2 \leq \sigma^2.$$

We have:

$$\begin{aligned}\mathbb{E}[f(x_{s+1}) - f(x_s)] &\leq -\eta \|\nabla f(x_s)\|^2 + \frac{\eta^2 l}{2} (\sigma^2 + \|\nabla f(x_s)\|^2) \\ &= \left(\frac{\eta^2 l}{2} - \eta\right) \|\nabla f(x_s)\|^2 + \frac{\eta^2 l \sigma^2}{2}.\end{aligned}$$

So:

$$\mathbb{E}[f(x_{s+1}) - f(x_s)] \leq \left(\frac{\eta^2 l}{2} - \eta\right) \|\nabla f(x_s)\|^2 + \frac{\eta^2 l \sigma^2}{2}. \quad (1)$$

By Strongly Convexity:

$$\begin{aligned}
\mathbb{E}(f(x_s) - f(x^*)) &\leq \mathbb{E} \left[ \langle \nabla f(x_s), x_s - x^* \rangle - \frac{\alpha}{2} \|x_s - x^*\|^2 \right] \\
&= \mathbb{E} [\langle \mathbb{E}[g(x_s)], x_s - x^* \rangle] - \frac{\alpha}{2} \mathbb{E} [\|x_s - x^*\|^2] \\
&= \mathbb{E} \left[ \left\langle \frac{1}{\eta} (x_s - x_{s+1}), x_s - x^* \right\rangle \right] - \frac{\alpha}{2} \mathbb{E} [\|x_s - x^*\|^2] \\
&= \frac{1}{2\eta} \mathbb{E} [\|x_s - x_{s+1}\|^2 + \|x_s - x^*\|^2 - \|x_{s+1} - x^*\|^2] - \frac{\alpha}{2} \mathbb{E} [\|x_s - x^*\|^2] \\
&= \left( -\frac{\alpha}{2} + \frac{1}{2\eta} \right) \mathbb{E} [\|x_s - x^*\|^2] - \frac{1}{2\eta} \mathbb{E} [\|x_{s+1} - x^*\|^2] + \frac{\eta}{2} \mathbb{E} [\|g(x_s)\|^2] \\
&\leq \left( -\frac{\alpha}{2} + \frac{1}{2\eta} \right) \mathbb{E} [\|x_s - x^*\|^2] - \frac{1}{2\eta} \mathbb{E} [\|x_{s+1} - x^*\|^2] + \frac{\eta}{2} (\sigma^2 + \|\nabla f(x_s)\|^2).
\end{aligned}$$

So:

$$\mathbb{E}[f(x_s) - f(x^*)] \leq \left( -\frac{\alpha}{2} + \frac{1}{2\eta} \right) \mathbb{E} [\|x_s - x^*\|^2] - \frac{1}{2\eta} \mathbb{E} [\|x_{s+1} - x^*\|^2] + \frac{\eta}{2} (\sigma^2 + \|\nabla f(x_s)\|^2). \quad (2)$$

Add up (1) and (2):

$$\begin{aligned}
\mathbb{E}[f(x_{s+1}) - f(x^*)] &\leq \left( -\frac{\alpha}{2} + \frac{1}{2\eta} \right) \mathbb{E} [\|x_s - x^*\|^2] - \frac{1}{2\eta} \mathbb{E} [\|x_{s+1} - x^*\|^2] \\
&\quad + \left( \frac{\eta^2 l \sigma^2}{2} + \frac{\eta \sigma^2}{2} \right) + \left( \frac{\eta^2 l}{2} - \frac{\eta}{2} \right) \|\nabla f(x_s)\|^2 \\
&\leq \left( -\frac{\alpha}{2} + \frac{1}{2\eta} \right) \mathbb{E} [\|x_s - x^*\|^2] - \frac{1}{2\eta} \mathbb{E} [\|x_{s+1} - x^*\|^2] + \eta \sigma^2
\end{aligned}$$

Do some simplification and we gain the final ans:

$$\mathbb{E}\|x_{s+1} - x^*\|^2 \leq (1 - \eta\alpha) \mathbb{E}\|x_s - x^*\|^2 - 2\eta \mathbb{E}[f(x_{s+1}) - f(x^*)] + 2\eta^2 \sigma^2 \quad (3)$$

### 1.2 b

Do telescope of the inequality in (b):

$$\begin{aligned}
\mathbb{E} \sum_{s=2}^{t+1} \lambda_s (f(x_s) - f(x^*)) &\leq \sum_{s=2}^{t+1} \frac{\lambda_s}{2\eta} \mathbb{E} [(1 - \alpha\eta) \|x_{s-1} - x^*\|^2 - \|x_s - x^*\|^2 + 2\eta^2 \sigma^2] \\
&\quad \text{Plug in the value of } \lambda_s \\
&\leq \eta\sigma^2 + \frac{1}{2\eta \sum_{s=2}^{t+1} (1 - \eta\alpha)^{t+1-s}} \sum_{s=2}^{t+1} [(1 - \alpha\eta)^{t+2-s} \|x_{s-1} - x^*\|^2 - (1 - \alpha\eta)^{t+1-s} \|x_s - x^*\|^2] \\
&= \eta\sigma^2 + \frac{1}{2\eta \sum_{s=2}^{t+1} (1 - \eta\alpha)^{t+1-s}} [(1 - \alpha\eta)^t \|x_1 - x^*\|^2 - (1 - \eta\alpha)^0 \|x_{t+1} - x^*\|^2] \\
&\leq \eta\sigma^2 + \frac{1}{2\eta \sum_{s=2}^{t+1} (1 - \eta\alpha)^{t+1-s}} e^{-\eta\alpha t} \|x_1 - x^*\|^2 \\
&\leq \eta\sigma^2 + \frac{e^{-\eta\alpha t} \|x_1 - x^*\|^2}{2\eta}
\end{aligned}$$

We can easily know that  $\eta\alpha \in (0, 1)$ . The last step relax the denominator to 1.

We prove that:

$$\mathbb{E} \sum_{s=2}^{t+1} \lambda_s (f(x_s) - f(x^*)) \leq \eta\sigma^2 + \frac{e^{-\eta\alpha t} \|x_1 - x^*\|^2}{2\eta} \quad (4)$$

### 1.3 c

Use convexity? or use strong cvx? simply relaxation in parameter may not work  
Do the same operations just like in note 11, consider the dominator in small t and large t. Plug in different  $\eta$  in the two separate terms. But there still exists a strange constant, which doesn't cause too much trouble.

## 2 Catalyst Acceleration for Finite-sum Problems