
CHAPTER 6. STOCHASTIC GRADIENT DESCENT

Kun Yuan

April 11, 2024

1 Problem formulation

This chapter considers the following stochastic optimization problem

$$\min_{x \in \mathbb{R}^d} f(x) = \mathbb{E}_{\xi \sim \mathcal{D}}[F(x; \xi)] \quad (1)$$

where $\xi \sim \mathcal{D}$ denotes the random data sample and \mathcal{D} denotes the data distribution. Since \mathcal{D} is typically unknown in machine learning, the closed-form of $f(x)$ is also unknown.

Notation. We introduce the following notations:

- Let $x^* := \arg \min_{x \in \mathbb{R}^d} \{f(x)\}$ be the optimal solution to problem (1).
- Let $f^* := \min_{x \in \mathbb{R}^d} \{f(x)\}$ be the optimal function value.
- Let $\mathcal{F}_k = \{x_k, \xi_{k-1}, x_{k-1}, \dots, \xi_0\}$ be the filtration containing all historical variables at and before iteration k . Note that ξ_k does not belong to \mathcal{F}_k .

2 Stochastic gradient descent

Since $f(x)$ does not have a closed-form, we cannot access its gradient. However, since $F(x; \xi)$ is known, we can use $\nabla_x F(x; \xi)$ to approximate the true gradient $\nabla f(x)$. Throughout this lecture, we let $\nabla F(x; \xi) = \nabla_x F(x; \xi)$ for notation simplicity. Given any arbitrary initialization variable x_0 , stochastic gradient descent (SGD) iterates as follows

$$x_{k+1} = x_k - \gamma \nabla F(x_k; \xi_k), \quad \forall k = 0, 1, 2, \dots \quad (2)$$

where γ is the learning rate, and $\xi_k \sim \mathcal{D}$ is a random data sampled at iteration k . Since ξ_k is a random variable for any $k = 0, 1, \dots$, each variable x_k is also a random variable for $k = 1, 2, \dots$.

3 Convergence analysis

To facilitate convergence analysis, we introduce the following assumption:

Assumption 3.1. Given the filtration \mathcal{F}_k , we assume

$$\mathbb{E}[\nabla F(x_k; \xi_k) | \mathcal{F}_k] = \nabla f(x_k) \quad (3)$$

$$\mathbb{E}[\|\nabla F(x_k; \xi_k) - \nabla f(x_k)\|^2 | \mathcal{F}_k] \leq \sigma^2 \quad (4)$$

The above assumption indicates that, conditioned on the filtration \mathcal{F}_k , the stochastic gradient $\nabla F(x_k; \xi_k)$ is an unbiased estimate on $\nabla f(x_k)$, and the variance is bounded by σ^2 . Under the above assumption, it is easy to verify that

$$\begin{aligned} \mathbb{E}[\|\nabla F(x_k; \xi_k)\|^2 | \mathcal{F}_k] &= \mathbb{E}[\|\nabla F(x_k; \xi_k) - \nabla f(x_k) + \nabla f(x_k)\|^2 | \mathcal{F}_k] \\ &= \|\nabla f(x_k)\|^2 + \mathbb{E}[\|\nabla F(x_k; \xi_k) - \nabla f(x_k)\|^2 | \mathcal{F}_k] \\ &\leq \|\nabla f(x_k)\|^2 + \sigma^2 \end{aligned} \quad (5)$$

where the second equality holds due to (3) and the last inequality holds due to (4).

3.1 Smooth and non-convex problem

Theorem 3.2. Suppose $f(x)$ is L -smooth and Assumption 3.1 holds. If $\gamma \leq 1/L$, SGD will converge at the following rate

$$\frac{1}{K+1} \sum_{k=0}^K \mathbb{E}[\|\nabla f(x_k)\|^2] \leq \frac{2\Delta_0}{\gamma(K+1)} + \gamma L \sigma^2, \quad (6)$$

where $\Delta_0 = f(x_0) - f^*$. If we further choose $\gamma = \left[\left(\frac{2\Delta_0}{(K+1)L\sigma^2} \right)^{-\frac{1}{2}} + L \right]^{-1}$, SGD converges as

$$\frac{1}{K+1} \sum_{k=0}^K \mathbb{E}[\|\nabla f(x_k)\|^2] \leq \sqrt{\frac{8L\Delta_0\sigma^2}{K+1}} + \frac{2L\Delta_0}{K+1}. \quad (7)$$

Remark. If $\sigma^2 = 0$, the stochastic gradient reduces to the true gradient, and hence, SGD reduces to GD. Substituting $\sigma^2 = 0$ to SGD rate (7), we recover the rate $O(L/K)$ for GD. In other words, our convergence rate for SGD is consistent with GD.

Proof. Since $f(x)$ is L -smooth, we have

$$\begin{aligned}
\mathbb{E}[f(x_{k+1})|\mathcal{F}_k] &\leq f(x_k) + \mathbb{E}[\langle \nabla f(x_k), x_{k+1} - x_k \rangle | \mathcal{F}_k] + \frac{L}{2} \mathbb{E}[\|x_{k+1} - x_k\|^2 | \mathcal{F}_k] \\
&= f(x_k) - \gamma \mathbb{E}[\langle \nabla f(x_k), \nabla F(x_k; \xi_k) \rangle | \mathcal{F}_k] + \frac{L\gamma^2}{2} \mathbb{E}[\|\nabla F(x_k; \xi_k)\|^2 | \mathcal{F}_k] \\
&\stackrel{(a)}{\leq} f(x_k) - \gamma(1 - \frac{L\gamma}{2}) \|\nabla f(x_k)\|^2 + \frac{L\gamma^2\sigma^2}{2} \\
&\stackrel{(b)}{\leq} f(x_k) - \frac{\gamma}{2} \|\nabla f(x_k)\|^2 + \frac{L\gamma^2\sigma^2}{2}
\end{aligned} \tag{8}$$

where inequality (a) holds due to Assumption 3.1, and inequality (b) holds if $\gamma \leq 1/L$. By taking expectations over the filtration \mathcal{F}_k , we have

$$\mathbb{E}[f(x_{k+1})] \leq \mathbb{E}[f(x_k)] - \frac{\gamma}{2} \mathbb{E}[\|\nabla f(x_k)\|^2] + \frac{L\gamma^2\sigma^2}{2} \tag{9}$$

which is equivalent to

$$\mathbb{E}[\|\nabla f(x_k)\|^2] \leq \frac{2}{\gamma} (\mathbb{E}[f(x_k)] - \mathbb{E}[f(x_{k+1})]) + \gamma L\sigma^2 \tag{10}$$

Taking averaging over $k = 0, 1, \dots, K$, we have

$$\frac{1}{K+1} \sum_{k=0}^K \mathbb{E}[\|\nabla f(x_k)\|^2] \leq \frac{2(f(x_0) - f^*)}{\gamma(K+1)} + \gamma L\sigma^2. \tag{11}$$

Defining $\Delta_0 := f(x_0) - f^*$, if we set

$$\gamma = \left[\left(\frac{2\Delta_0}{(K+1)L\sigma^2} \right)^{-\frac{1}{2}} + L \right]^{-1}, \tag{12}$$

it then holds that

$$\gamma \leq \min \left\{ \frac{1}{L}, \gamma_1 \right\}, \quad \text{where} \quad \gamma_1 = \left(\frac{2\Delta_0}{(K+1)L\sigma^2} \right)^{\frac{1}{2}}. \tag{13}$$

Substituting (12) and (13) into (11), we have

$$\begin{aligned}
\frac{1}{K+1} \sum_{k=0}^K \mathbb{E}[\|\nabla f(x_k)\|^2] &\leq \frac{2(f(x_0) - f^*)}{\gamma(K+1)} + \gamma_1 L\sigma^2 \\
&= 2\sqrt{\frac{2L\Delta_0\sigma^2}{K+1}} + \frac{2L\Delta_0}{K+1}.
\end{aligned} \tag{14}$$

□