

# Stochastic Gradient Descent

$$f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x) \quad (\text{e.g. } x \in \mathbb{R}^n \text{ weights in NN})$$

Gradient Descent  $x_{k+1} = x_k - \alpha \nabla f(x_k)$ .  $O(n)$

Stochastic Gradient Descent (SGD):  $x_{k+1} = x_k - \alpha_k \nabla f_{i_k}(x_k)$   $O(1)$   
 $i_k \in \{1, \dots, n\}$  chosen at random.  
 $\uparrow$  uniformly

Goal: Convergence theory for SGD.

Define  $\eta_k = \nabla f_{i_k}(x_k) - \nabla f(x_k)$ .

Then SGD is

$$x_{k+1} = x_k - \alpha_k (\nabla f(x_k) + \eta_k)$$

$\downarrow$  noise

Conditional expectation

$$\mathbb{E}_k(\nabla f_{i_k}(x_k)) = \frac{1}{n} \sum_{i=1}^n \nabla f_i(x_k) = \nabla f(x_k)$$

Since  $\mathcal{M}_k = \nabla f_{i_k}(x_k) - \nabla f(x_k)$ ,

$$\mathbb{E}_k(\mathcal{M}_k) = \mathbb{E}_k(\nabla f_{i_k}(x_k) - \nabla f(x_k)) = 0$$

Assumption (\*)  $\mathbb{E}_k(\|\mathcal{M}_k\|^2) \leq \sigma^2$ ,  $\sigma \geq 0$   
on noise variance

$$\begin{aligned}\mathbb{E}_k(\|\mathcal{M}_k\|^2) &= \mathbb{E}_k(\|\nabla f_{i_k}(x_k) - \nabla f(x_k)\|^2) \\ &= \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x_k) - \nabla f(x_k)\|^2 \\ &= \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x_k)\|^2 - \frac{2}{n} \sum_{i=1}^n \nabla f(x_k)^T \nabla f_i(x_k) \\ &\quad + \frac{1}{n} \sum_{i=1}^n \|\nabla f(x_k)\|^2 \\ &= \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x_k)\|^2 - 2 \nabla f(x_k)^T \underbrace{\left[ \frac{1}{n} \sum_{i=1}^n \nabla f_i(x_k) \right]}_{\nabla f(x_k)} \\ &\quad + \|\nabla f(x_k)\|^2 \\ &= \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x_k)\|^2 - \|\nabla f(x_k)\|^2\end{aligned}$$

(\*) is equivalent to  $\forall x \in \mathbb{R}^n$

$$\frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x)\|^2 \leq \|\nabla f(x)\|^2 + \sigma^2$$

Assume  $f$  is  $L$ -Lipschitz, so that

$$f(x_{k+1}) \leq f(x_k) + \nabla f(x_k)^T (x_{k+1} - x_k) + \frac{L}{2} \|x_{k+1} - x_k\|^2$$

$$\text{For SGD, } x_{k+1} - x_k = -\alpha_k \nabla f_{i_k}(x_k)$$

$$f(x_{k+1}) \leq f(x_k) - \alpha_k \nabla f(x_k)^T \nabla f_{i_k}(x_k) + \frac{L\alpha_k^2}{2} \|\nabla f_{i_k}(x_k)\|^2$$

Take conditional expectation on both sides

$$\mathbb{E}_k(f(x_{k+1})) \leq f(x_k) - \alpha_k \nabla f(x_k)^T \underbrace{\mathbb{E}_k[\nabla f_{i_k}(x_k)]}_{\nabla f(x_k)} + \frac{L\alpha_k^2}{2} \underbrace{\mathbb{E}_k[\|\nabla f_{i_k}(x_k)\|^2]}$$

$$\hookrightarrow = \frac{1}{n} \sum_{i=1}^{\hat{n}} \|\nabla f_i(x_k)\|^2 \leq \|\nabla f(x_k)\|^2 + \sigma^2$$

by  $(\phi)$

$$\mathbb{E}_k(f(x_{k+1})) \leq f(x_k) - \alpha_k \|\nabla f(x_k)\|^2 + \frac{L\alpha_k^2}{2} \|\nabla f(x_k)\|^2 + \frac{L\alpha_k^2 \sigma^2}{2}$$

Assume  $\alpha_k \leq \frac{1}{L}$

then  $\frac{L\alpha_k^2}{2} \leq \frac{L}{2} \cdot \frac{1}{L} \cdot \alpha_k = \frac{\alpha_k}{2}$

Then

$$\mathbb{E}_k(f(x_{k+1})) \leq f(x_k) - \frac{\alpha_k}{2} \|\nabla f(x_k)\|^2 + \frac{L\alpha_k^2 \sigma^2}{2}$$

Rearrange:

↑  
Noise in  
SGD

$$\frac{\alpha_k}{2} \|\nabla f(x_k)\|^2 \leq f(x_k) - \underbrace{\mathbb{E}_k(f(x_{k+1}))}_{\text{only over } i_k} + \frac{L\alpha_k^2 \sigma^2}{2}$$

Take expectation over everything  $i_0, i_1, \dots, i_{k-1}$   
and use iterated expectation

$$\mathbb{E} \mathbb{E}_k (f(x_{k+1})) = \mathbb{E} (f(x_{k+1}))$$

$$\frac{\alpha_k}{2} \mathbb{E} \|\nabla f(x_k)\|^2 \leq \mathbb{E}(f(x_k)) - \mathbb{E}(f(x_{k+1})) + \frac{L\alpha_k^2\sigma^2}{2}$$

Sum both sides from  $k=0, \dots, t-1$

$$\frac{1}{2} \sum_{k=0}^{t-1} \alpha_k \mathbb{E} \|\nabla f(x_k)\|^2 \leq \underbrace{\sum_{k=0}^{t-1} (\mathbb{E}(f(x_k)) - \mathbb{E}(f(x_{k+1})))}_{\text{Telescoping}} + \frac{L\sigma^2}{2} \sum_{k=0}^{t-1} \alpha_k^2$$

$$\begin{aligned} \frac{1}{2} \sum_{k=0}^{t-1} \alpha_k \mathbb{E} \|\nabla f(x_k)\|^2 &\leq f(x_0) - \mathbb{E}(f(x_t)) + \frac{L\sigma^2}{2} \sum_{k=0}^{t-1} \alpha_k^2 \\ &\leq f(x_0) - f_* + \frac{L\sigma^2}{2} \sum_{k=0}^{t-1} \alpha_k^2 \end{aligned}$$

When  $f_* = \min_{\mathcal{R}^n} f$ . Now use that

$$\frac{\sum_{k=0}^{t-1} \alpha_k \mathbb{E} \|\nabla f(x_k)\|^2}{\sum_{k=0}^{t-1} \alpha_k} \geq \min_{0 \leq k \leq t-1} \mathbb{E} \|\nabla f(x_k)\|^2$$

weighted average

Here

$$\sum_{k=0}^{t-1} \alpha_k \mathbb{E} \|\nabla f(x_k)\|^2 \geq \left( \min_{0 \leq k \leq t-1} \mathbb{E} \|\nabla f(x_k)\|^2 \right) \sum_{k=0}^{t-1} \alpha_k$$

Plus this in above to get

$$\left( \min_{0 \leq k \leq t-1} \mathbb{E} \|\nabla f(x_k)\|^2 \right) \sum_{k=0}^{t-1} \alpha_k \leq 2 \left( f(x_0) - f_* + \frac{L\sigma^2}{2} \sum_{k=0}^{t-1} \alpha_k^2 \right)$$

$$\min_{0 \leq k \leq t-1} \mathbb{E} \|\nabla f(x_k)\|^2 \leq \frac{2(f(x_0) - f_*)}{\sum_{k=0}^{t-1} \alpha_k} + \frac{L\sigma^2}{2} \frac{\sum_{k=0}^{t-1} \alpha_k^2}{\sum_{k=0}^{t-1} \alpha_k}$$

Remarks

① Constant step size  $\alpha_k = \alpha$  SGD Noise

$$\min_{0 \leq k \leq t-1} \mathbb{E} \|\nabla f(x_k)\|^2 \leq \frac{2(f(x_0) - f_*)}{\alpha t} + \frac{L\sigma^2 \alpha t}{2 \alpha t}$$

Noise does not vanish as  $t \rightarrow \infty$ . SGD with

$$\frac{L\sigma^2 \alpha}{2}$$

constant step size does not converge.

$$\textcircled{2} \quad \alpha_k = \frac{\alpha}{k+1}, \quad \sum_{k=0}^{t-1} \alpha_k = O(\alpha \log(t))$$

Compare  $\sum_{k=0}^{t-1} \frac{1}{k+1}$   
to  $\int_0^t \frac{1}{x+1} dx$

$$\sum_{k=0}^{t-1} \alpha_k^2 = \alpha^2 \sum_{k=0}^{t-1} \frac{1}{(k+1)^2} \leq \frac{\pi^2}{6} \alpha^2$$

Here error term is

$$\frac{L\alpha^2}{2} \frac{\sum_{k=0}^{t-1} \alpha_k^2}{\sum_{k=0}^{t-1} \alpha_k} = \frac{L\alpha^2}{2} O\left(\frac{1}{\log(t)}\right) \xrightarrow{t \rightarrow \infty} 0$$

$$\textcircled{3} \quad \alpha_k = \frac{\alpha}{\sqrt{k+1}}, \quad \sum_{k=0}^{t-1} \alpha_k = \alpha \sum_{k=0}^{t-1} \frac{1}{\sqrt{k+1}}$$
$$\sum_{k=0}^{t-1} \alpha_k^2 = \alpha^2 \sum_{k=0}^{t-1} \frac{1}{k+1} \sim \alpha \int_0^t \frac{1}{\sqrt{x+1}} dx$$
$$= \alpha O(\log(t)) \sim O(\alpha \sqrt{t})$$

Error term

$$\frac{L\sigma^2}{2} \frac{\sum_{k=0}^{t-1} \alpha_k^2}{\sum_{k=0}^{t-1} \alpha_k} = \frac{L\sigma^2}{2} O\left(\frac{\log(t)}{\sqrt{t}}\right)$$

much better

Note: worse than  $O\left(\frac{1}{t}\right)$  when  $\sigma=0$ .

---

## Strong Convexity

Assume  $f$  is  $\mu$ -strongly convex,  
so the PL-inequality

$$\frac{1}{2} \|\nabla f(x)\|^2 \geq \mu (f(x) - f_*)$$

Recall:

$$\mathbb{E}_k (f(x_{k+1})) \leq f(x_k) - \frac{\alpha_k}{2} \|\nabla f(x_k)\|^2 + \frac{L\alpha_k^2 \sigma^2}{2}$$

$$PL \leq f(x_k) - \alpha_k \mu (f(x_k) - f_*) + \frac{L\alpha_k^2 \sigma^2}{2}$$

Subtract  $f_*$  from both sides



$$\mathbb{E}_k (f(x_{k+1}) - f_*) \leq (1 - \alpha_k \mu) (f(x_k) - f_*) + \frac{L \alpha_k^2 \sigma^2}{2}$$

Take expectations on both sides over all  $i_0, i_1, i_2, \dots, i_{k-1}$  and set

$$e_k = \mathbb{E} (f(x_k) - f_*)$$

$$e_{k+1} \leq (1 - \alpha_k \mu) e_k + \frac{L \alpha_k^2 \sigma^2}{2}$$

Choose  $\alpha_k = \frac{1}{\mu(k+2)}$ ,  $\alpha_k \mu = \frac{1}{k+2}$

$$e_{k+1} \leq \left(1 - \frac{1}{k+2}\right) e_k + \frac{L \sigma^2}{2(k+2)^2 \mu^2}$$

Claim:  $e_k = o\left(\frac{1}{k}\right)$

let's consider  $f_k = \frac{\mu^2}{\frac{1}{2} L \sigma^2} e_{k-2}$

$$\left\{ \begin{array}{l} f_{k+1} \leq (1 - \frac{1}{k}) f_k + \frac{1}{k^2}, \quad k \geq 2 \\ f_2 = \frac{\mu^2}{\frac{1}{2} L \sigma^2} e_0 \end{array} \right.$$

Claim:  $f_k \leq \frac{C}{k}$ .

Proof:  $f_{k+1} \leq (\frac{k-1}{k}) f_k + \frac{1}{k^2}$

$$\leq (\frac{k-1}{k}) \left( (\frac{k-2}{k-1}) f_{k-1} + \frac{1}{(k-1)^2} \right) + \frac{1}{k^2}$$

$$= (\frac{k-2}{k}) f_{k-1} + \frac{1}{k} \left( \frac{1}{k-1} + \frac{1}{k} \right)$$

$$\leq (\frac{k-2}{k}) \left( (\frac{k-3}{k-2}) f_{k-2} + \frac{1}{(k-2)^2} \right) + \frac{1}{k} \left( \frac{1}{k-1} + \frac{1}{k} \right)$$

$$\begin{aligned}
&= \left(\frac{k-3}{k}\right) f_{k-2} + \frac{1}{k} \left(\frac{1}{k-2} + \frac{1}{k-1} + \frac{1}{k}\right) \\
&\leq \left(\frac{1}{k}\right) f_2 + \frac{1}{k} \left(\frac{1}{2} + \frac{1}{3} + \dots + \frac{1}{k}\right) \\
&\leq \frac{1}{k} \frac{\mu^2}{\frac{1}{2} \log^2} e_0 + \frac{1}{k} \underbrace{\sum_{j=2}^k \frac{1}{j}}_{\log(k)} \\
&= O\left(\frac{\log(k)}{k}\right).
\end{aligned}$$

Main Point: Compare to  $O\left(\frac{1}{\sqrt{k}}\right)$  rate for general functions.

linear rate.

Recall rate is  $O(r^k)$   
for gradient descent  
on  $\mu$ -strongly convex functions  
with  $0 < r < 1$ .

SGD does not see  
strong convexity.