

Problem Class :  $f: \mathbb{R}^n \rightarrow \mathbb{R}$

$$f^* = \inf_{x \in \mathbb{R}^n} f(x) > -\infty$$

Goal: Find  $\bar{x}$ :  $\|f'(\bar{x})\|_* \leq \varepsilon$ .

Smoothness:  $\|f'(y) - f'(x)\|_* \leq L \cdot \|y - x\| \quad \forall x, y \in \mathbb{R}^n$ .

Init.  $x_0 \in \mathbb{R}^n$ . Iterate,  $k \geq 0$ :

$$x_{u+1} = x_k^+(M_h) = \underset{y \in \mathbb{R}^n}{\operatorname{argmin}} \left[ f(x_k) + \langle f'(x_k), y - x_k \rangle + \frac{M_h}{2} \|y - x_k\|^2 \right]$$

Norm is Euclidean:  $x_{u+1} = x_u - \frac{1}{M} f'(x_u)$ .

For  $M \geq L$ :  $f(x_u) - f(x_u^+(M)) \geq \frac{1}{2M} \|f'(x_u)\|_*^2$ .

In theory: M := L.

## Adaptive Search.

$$M_k \approx L$$


---

## Gradient Method with Adaptive Search.

Init.:  $x_0 \in \mathbb{R}^n$ ,  $\epsilon > 0$ ,  $M_0 > 0$

Iterations,  $k \geq 0$ :

1. If  $\|f'(x_k)\|_* \leq \epsilon$  then return  $x_k$ .

2. For  $t \geq 0$  iterate:

- Set  $M_k^+ = M_k \cdot 2^t$ .
- Try gradient step:  $x_k^+ := x_k + M_k^+ f'(x_k)$
- If  $f(x_k) - f(x_k^+) \geq \frac{1}{2M_k^+} \|f'(x_k)\|_*^2$  then break and go to step 3.

3.  $x_{k+1} = x_k^+$ ,  $M_{k+1} = M_k^+ \cdot \frac{1}{2}$ .

---

As soon as  $M_k \cdot 2^t \geq L$  then we exit from inner loop.

Proposition I For  $k \geq 0$ :

$$M_k \leq M_* := \max\{M_0, L\}.$$

Proof For  $k=0$ .  $M_0 \leq M_*$ . true.

Consider  $k \geq 0$ .

Denote  $t_k \geq 0$  the value of parameter  $t$  on step 2 that triggered the break.

$$M_{k+1} = M_k \cdot 2^{t_k - 1}.$$

Case 1  $t_k = 0$ .  $\Rightarrow M_{k+1} = \frac{1}{2} M_k \leq \frac{1}{2} M^* \leq M^*$ .

$t_k > 0$ .  $\Rightarrow M_{k+1} < L \leq M^*$ .  $\square$

Corollary: To find  $\bar{x}$ :  $(f'(\bar{x}))_{L^*} \leq \varepsilon$  it's enough

to do

$$K = \left\lceil \frac{4 \cdot \max\{M_0, L\} (f(x_0) - f^*)}{\varepsilon^2} \right\rceil.$$

Proposition 2 let  $N_K$  the total number of oracle calls during the first  $K$  iterations.

$$N_K \leq 2K + \max\{0, 1 + \log_2 \frac{L}{M_0}\}.$$

Proof.

$$2^{t_k - 1} = \frac{M_{k+1}}{M_k} \Rightarrow t_k = 1 + \log_2 \frac{M_{k+1}}{M_k}.$$

$$\begin{aligned} N_K &= \sum_{i=0}^{K-1} (1 + t_i) = 2K + \sum_{i=0}^{K-1} \log_2 \frac{M_{i+1}}{M_i} \\ &= 2K + \log_2 \frac{M_K}{M_0} \leq 2K + \log_2 \frac{M^*}{M_0}. \end{aligned}$$

$\square$ .

## Stochastic Gradient Method.

### Example (Expensive)

$$f(x) = \frac{1}{N} \sum_{i=1}^N f_i(x), \quad f_i : \mathbb{R}^n \rightarrow \mathbb{R}$$

$N$ - size of dataset.

$$f'(x) = \frac{1}{N} \sum_{i=1}^N f'_i(x) \quad - \text{expensive to spend } O(N)$$

but, if it's affordable

$$f'_i(x)$$

### Example (Impossible)

$$f(x) = \mathbb{E}_{\zeta} F(x, \zeta)$$

$F(\cdot, \zeta)$  is smooth, for every  $\zeta$ .

### Example (We can, but we don't want)

e.g. when  $f_i(x)$  reveals user data.

## Stochastic Oracle

$\| \cdot \| = \| \cdot \|_2$  Euclidean.

Stochastic first-order oracle by

$$g(x, \xi) \in \mathbb{R}^n$$

Ex. Finite sum:  $\xi = i \sim \{1, \dots, N\}$

$$g(x, \xi) = f_i(x)$$

1. Unbiased estimator:  $\forall x \in \mathbb{R}^n$

$$\mathbb{E}_{\xi} [g(x, \xi)] = f(x)$$

2. Bounded variance:  $\sigma > 0$  s.t.  $\forall x \in \mathbb{R}^n$

$$\mathbb{E}_{\xi} \|g(x, \xi) - f(x)\|^2 \leq \sigma^2 \quad (2)$$

$$\|g(x, \xi) - f(x)\|^2 = \|g(x, \xi)\|^2 + \|f(x)\|^2 - 2\langle g(x, \xi), f(x) \rangle$$

$$\sigma^2 \geq \mathbb{E} \|g(x, \xi) - f'(x)\|^2 = \mathbb{E} \|g(x, \xi)\|^2 - \mathbb{E} f'(x))^2$$

$$(2) \Leftrightarrow \mathbb{E} \|g(x, \xi)\|^2 \leq \|f'(x)\|^2 + \sigma^2.$$

### Stochastic Gradient Method.

Initialization:  $x_0 \in \mathbb{R}^n$ ,  $M > 0$ , number of iterations  
 $K \geq 1$ .

For  $k = 0 \dots K-1$ :

- 1. Sample  $\xi_k$
- 2. Compute  $g_k = g(x_k, \xi_k)$
- 3. Gradient Step:

$$x_{k+1} = x_k - \frac{1}{M} g_k.$$

Sample  $j \in \{0, \dots, K-1\}$  uniformly  
 and return  $\bar{x}_K = x_j$ .

Problem:

$$\min_x f(x)$$

Find a point  $\bar{x}$ :  $\mathbb{E} \|f'(\bar{x})\| < \epsilon$ .

Proposition let  $M > 0$ .

$$x_{un} = x_u - \frac{1}{M} g_u \text{ with } g_u = g(x_u, \xi_u)$$

Then

$$\mathbb{E}_{\xi_u} [f(x_u) - f(x_{un})] \geq \frac{1}{M} \|f'(x_u)\|^2 \cdot \underbrace{\left(1 - \frac{L}{2M}\right)}_{-\frac{L}{2M^2} \sigma^2}.$$

Remark  $M \geq L \Leftrightarrow 1 - \frac{L}{2M} \geq \frac{1}{2}$ .

We obtain:

$$\mathbb{E}_{\xi_u} [f(x_u) - f(x_{un})] \geq \frac{1}{2M} \|f'(x_u)\|^2 - \frac{L}{2M^2} \sigma^2$$

If we  $\sigma = 0 \Rightarrow$  the same as for  $\mathcal{P}$ .

Stochastic Gradient Descent (SGD)

Proof Smoothness:

$$\begin{aligned} f(x_{un}) &\leq f(x_u) + \underbrace{\langle f'(x_u), x_{un} - x_u \rangle}_{= -\frac{1}{M} g_u} + \frac{L}{2} \|x_{un} - x_u\|^2 \\ &= f(x_u) - \frac{1}{M} g_u \end{aligned}$$

$$= f(x_u) - \frac{1}{M} \langle f'(x_u), g_u \rangle + \frac{L}{2M^2} \|g_u\|^2$$

Rearrange:

$$\mathbb{E}_{\tilde{\gamma}_u} [f(x_u) - f(\bar{x}_u)] \geq \mathbb{E}_{\tilde{\gamma}} \left[ \frac{1}{M} \langle f'(\bar{x}_u), g_u \rangle - \frac{L}{2M^2} \|g_u\|^2 \right]$$

$$= \frac{1}{M} \|f'(\bar{x}_u)\|^2 - \frac{L}{2M^2} \|f'(\bar{x}_u)\|^2 - \frac{L}{2M^2} \sigma^2 \quad \square.$$

Theorem  $M \geq L$ . Consider  $K \geq 1$  iterations.

$$\mathbb{E} \|f'(\bar{x}_K)\|^2 \leq \frac{2M(f(x_0) - f^*)}{K} + \frac{L}{M} \sigma^2.$$

Proof:

$\mathbb{E}[\cdot]$  the expectation w.r.t. all randomness

$\tilde{\gamma}_1, \dots, \tilde{\gamma}_K, j$ .

$\mathbb{E}_{\tilde{\gamma}} [\cdot]$  the expectation w.r.t all

$\tilde{\gamma}_1, \dots, \tilde{\gamma}_K$ .

$$\mathbb{E}_{\mathcal{S}}[f(x_u) - f(x_{u+1})] \geq \frac{1}{2M} \mathbb{E}_{\mathcal{S}}[\|f'(x_u)\|^2] - \frac{L}{2M^2} \sigma^2.$$

Telescope it for  $K \geq 1$  iterations:

$$f(x_0) - f^* \geq \mathbb{E}_{\mathcal{S}}[f(x_0) - f(x_K)] \geq \frac{1}{2M} \sum_{i=0}^{K-1} \mathbb{E}_{\mathcal{S}}[\|f'(x_i)\|^2] - K \cdot \frac{L}{2M^2} \sigma^2.$$

$$\frac{2M}{K} (f(x_0) - f^*) \geq \frac{1}{K} \sum_{i=0}^{K-1} \mathbb{E}_{\mathcal{S}}[\|f'(x_i)\|^2]$$

$$- \frac{L}{M} \sigma^2 =$$

$$= \mathbb{E}[\|f'(\bar{x}_K)\|^2] - \frac{L}{M} \sigma^2 \quad \square$$

## Stepsize Tuning

$$M \geq L.$$

$$\mathbb{E}[\|f'(\bar{x}_n)\|^2] \leq \frac{2MF_0}{K} + \frac{L}{M}\sigma^2$$

where  $F_0 = f(x_0) - f^*$ .

Naive :  $M = L$ :

$$\mathbb{E}[\|f'(\bar{x}_n)\|^2] \leq \frac{2LF_0}{K} + \sigma^2.$$

Now:

$$1. \quad \frac{2MF_0}{K} \leq \frac{\epsilon^2}{2}$$

$$2. \quad \frac{L}{M}\sigma^2 \leq \frac{\epsilon^2}{2} \Rightarrow M = L \cdot \max\left\{1, \frac{2\sigma^2}{\epsilon^2}\right\}$$

$$(1) \quad K = 1 + \left\lfloor \frac{4MF_0}{\epsilon^2} \right\rfloor \leq 1 + \frac{2\sigma^2}{\epsilon^2}$$

Corollary To find a random point  $\bar{x} \in \mathbb{R}^n$

s.t.  $\mathbb{E}\|f'(\bar{x})\| \leq \epsilon$  it's enough  
to do

$$K = O(L \cdot (f(x_0) - f^*) \cdot \left[ \frac{1}{\varepsilon^2} + \frac{\sigma^2}{\varepsilon^4} \right])$$

stochastic oracle calls.

⇒ optimal complexity.