

## Summary

$$\min_{x \in Q} f(x)$$

$$Q \subseteq \mathbb{R}^n \text{ - bounded, } D = \text{diam}(Q)$$

$$\|f'(x)\| \leq M \quad \forall x \in Q$$

## Stochastic Subgradient Method

$$x_{k+1} = \pi_Q \left( x_k - \frac{1}{\beta_k} \underline{g}_k \right), \quad \underline{g}_k \text{ - stochastic } \overset{\text{(sub)}}{\text{gradient}}$$

$$\beta_k := \frac{1}{D} \sqrt{\|g_{01}\|^2 + \dots + \|g_{0k}\|^2}$$

$$\mathbb{E}[f(\bar{x}_k) - f^*] \leq \frac{(\sigma + M)D}{\sqrt{k}}$$

$\sigma$ -variance:  $\mathbb{E} \|g_k - f'(x_k)\|^2 \leq \underline{\sigma^2}$

"Stochasticity" & "Non-smoothness" (convex)

$f: \mathbb{R}^n \rightarrow \mathbb{R}$  convex and smooth:

$$\|f'(y) - f'(x)\| \leq L \|y - x\|$$

unconstrained minimization:

$$f'(x^*) = 0$$

## Simple Analysis

$$x := x^*$$

$$\|f'(y)\| \leq L \|y - x^*\| \leq LD = M$$

$$D \geq \|x_0 - x^*\|$$

Extra constraint:  $Q = \{x \in \mathbb{R}^n : \|x - x_0\| \leq D\}$

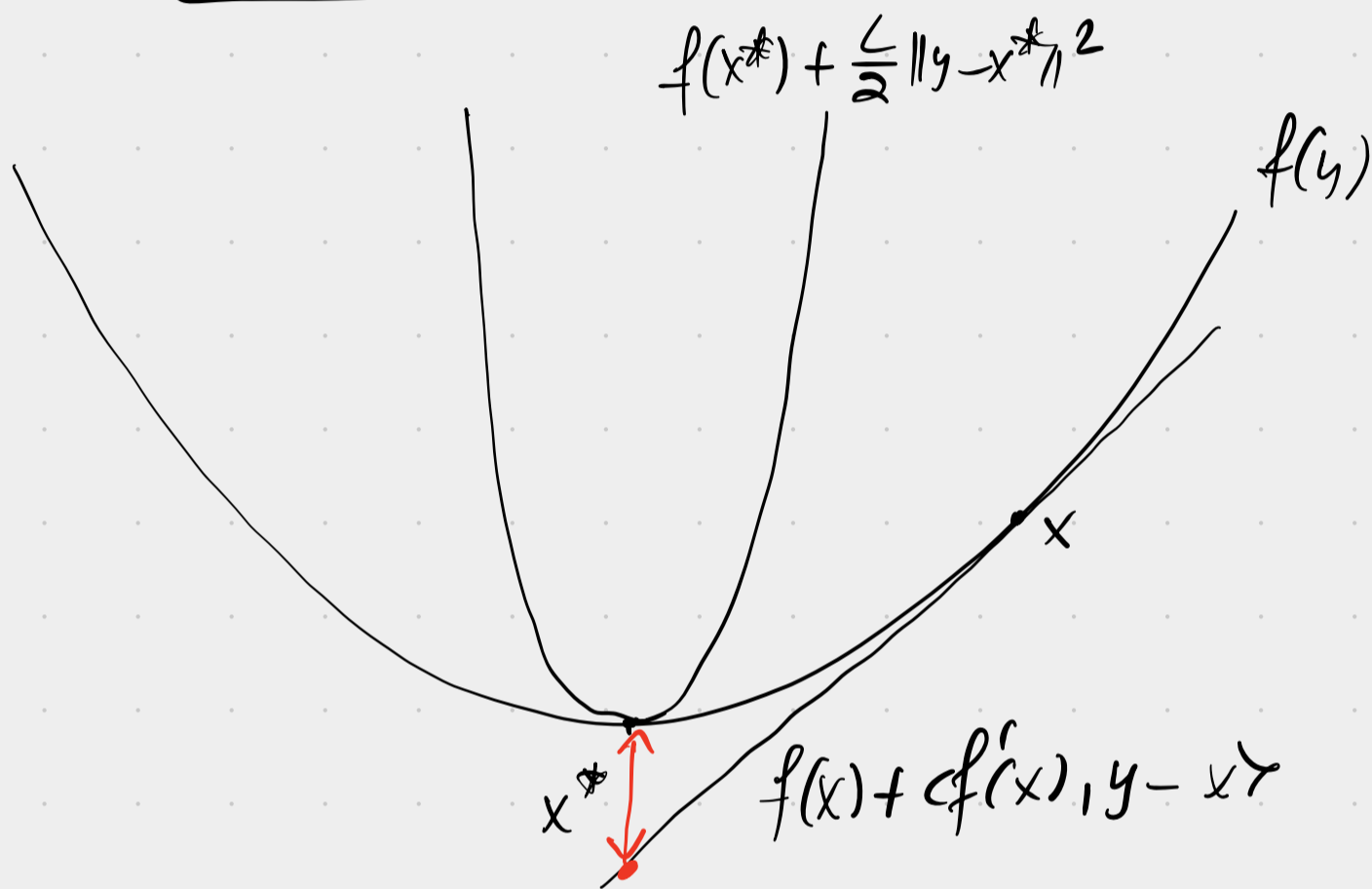
# Smooth Problems

Useful Inequality:  $\forall x, y \in \mathbb{R}^n$ :

$$f(y) \geq f(x) + \langle f'(x), y-x \rangle + \frac{1}{2L} \|f'(y) - f'(x)\|^2$$

$y := x^*$

$$f(x^*) \geq f(x) + \langle f'(x), x^* - x \rangle + \frac{1}{2L} \|f'(x)\|^2 \quad \forall x$$



$$f(x) + \langle f'(x), y-x \rangle \leq f(y) \leq f(x^*) + \frac{L}{2} \|x^* - y\|^2 \quad \forall y$$

$$\star f(x^*) \geq f(x) + \langle f'(x), x^* - x \rangle + \underbrace{\langle f'(x), y - x^* \rangle}_{\rightarrow \max} - \frac{L}{2} \|y - x^*\|^2$$

$= \frac{1}{2L} \|f'(x)\|^2$

Rotate: Fix  $x \in \mathbb{R}^n$ ; define

$$\varphi(y) = f(y) - \langle f'(x), y \rangle \Rightarrow \varphi'(y) = f'(y) - f'(x)$$

$$x_\varphi^* := x, \quad \varphi^* = f(x) - \langle f'(x), x \rangle$$

$$\underline{f(x) - \langle f'(x), x \rangle} \geq \underline{f(y) - \langle f'(x), y \rangle} + \langle f'(y) - f'(x), x - y \rangle$$

$$+ \frac{1}{2L} \|f'(y) - f'(x)\|^2.$$

Corollary  $\forall x, y \in \mathbb{R}^n$ ,  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ , convex and smooth.

$$f(y) \geq f(x) + \langle f'(x), y-x \rangle + \frac{1}{2L} \|f'(y) - f'(x)\|^2.$$

## Stochastic Gradient Method

$$x_{u+1} = x_u - \gamma g_u, \quad g_u = g(x_u, \xi_u)$$

$$\mathbb{E} g_u = f'(x_u)$$

$$\mathbb{E} \|g_u - f'(x_u)\|^2 \leq \sigma^2$$

$$\frac{1}{2} \|x_{u+1} - x^*\|^2 = \frac{1}{2} \|x_u - x^* - \gamma g_u\|^2 =$$

$$= \frac{1}{2} \|x_u - x^*\|^2 + \frac{\gamma^2}{2} \|g_u\|^2 - \gamma \langle g_u, x_u - x^* \rangle$$

Rearranging :

$$\frac{1}{2} \|x_u - x^*\|^2 + \frac{\gamma^2}{2} \|g_u\|^2 = \frac{1}{2} \|x_{u+1} - x^*\|^2 + \gamma \langle g_u, x_u - x^* \rangle$$

$$\mathbb{E}_{\xi_u} \|g_u - f'(x_u)\|^2 = \mathbb{E}_{\xi_u} \|g_u\|^2 - \|f'(x_u)\|^2$$

$$\Rightarrow \mathbb{E} \|g_u\|^2 = \mathbb{E} \|g_u - f'(x_u)\|^2 + \mathbb{E} \|f'(x_u)\|^2$$

$$\mathbb{E}_{\xi_u} \langle g_u, x_u - x^* \rangle = \langle f'(x_u), x_u - x^* \rangle \geq f(x_u) - f^* + \frac{1}{2L} \|f'(x_u)\|^2$$

$$\mathbb{E} \langle g_u, x_u - x^* \rangle \geq \mathbb{E} [f(x_u) - f^*] + \frac{1}{2L} \mathbb{E} \|f'(x_u)\|^2$$

Denote:  $R_k^2 := \mathbb{E} \|x_k - x^*\|^2$

Lemma :

$$\begin{aligned} \frac{1}{2}R_n^2 - \frac{1}{2}R_{n+1}^2 + \underbrace{\frac{\gamma^2}{2} \mathbb{E} \|g_n - f'(x_n)\|^2}_{\text{variance}} + \underbrace{\frac{\gamma^2}{2} \mathbb{E} \|f'(x_n)\|^2}_{\text{smoothness}} &\geq \\ &\geq \gamma \mathbb{E}[f(x_n) - f^*] + \underbrace{\frac{\gamma}{2L} \mathbb{E} \|f'(x_n)\|^2}_{\text{smoothness}} \end{aligned}$$

Corollary  $\gamma \leq \frac{1}{L}$  :

$$\frac{1}{2}R_n^2 - \frac{1}{2}R_{n+1}^2 + \frac{\gamma^2}{2} \mathbb{E} \|g_n - f'(x_n)\|^2 \geq \gamma \mathbb{E}[f(x_n) - f^*].$$

Apply variance assumption:

$$\frac{1}{2}R_n^2 - \frac{1}{2}R_{n+1}^2 + \frac{\gamma^2}{2}\sigma^2 \geq \gamma \mathbb{E}[f(x_n) - f^*]$$

Telescope for the first  $k$  iterations.

$$\underbrace{\frac{1}{2}R_0^2 - \frac{1}{2}R_k^2}_{\wedge} + \frac{\gamma^2}{2}\sigma^2 k \geq \gamma \sum_{i=0}^{k-1} \mathbb{E}[f(x_i) - f^*]$$

$$\frac{1}{2}R_0^2$$

$$\geq \gamma k \mathbb{E}[f(\bar{x}_k) - f^*]$$

Jensen's inequality.

Finally:

$$\mathbb{E}[f(\bar{x}_k) - f^*] \leq \underbrace{\frac{R_0^2}{2\gamma k}}_{\leq \frac{\epsilon}{2}} + \underbrace{\frac{\gamma\sigma^2}{2}}_{\leq \frac{\epsilon}{2}} \rightarrow \min_{\gamma}$$

$$\frac{R_0^2}{2\gamma_*^2 k} = \frac{\sigma^2}{2} \Rightarrow \gamma_* = \frac{R_0}{\sigma\sqrt{k}}, \quad \gamma \leq \frac{1}{L}$$

$$f := \min \left\{ \frac{1}{L}, \frac{R_0}{\sigma\sqrt{k}} \right\}$$

$$\frac{R_0^2}{2k} \max \left\{ L, \frac{\sigma\sqrt{k}}{R_0} \right\} + \frac{R_0}{\sigma\sqrt{k}} \cdot \frac{\sigma^2}{2} = O \left( \frac{LR_0^2}{k} + \frac{\sigma R_0}{\sqrt{k}} \right)$$

Theorem:

$$\mathbb{E}[f(X_k) - f^*] \leq O \left( \frac{LR_0^2}{k} + \frac{\sigma R_0}{\sqrt{k}} \right) \leq \varepsilon$$

$$k = O \left( \frac{LR_0^2}{\varepsilon} + \left[ \frac{\sigma R_0}{\varepsilon} \right]^2 \right)$$

stochastic oracle calls.

## Variance Reduction

### 1. Mini-Batching.

$m$  stochastic samples

$$g_u = \frac{1}{m} \sum_{i=1}^m g(x_u, \xi_{k,i}) \Rightarrow \mathbb{E}_{\xi_u} g_u = f'(x_u)$$

$$\mathbb{E} \|g_u - f'(x_u)\|^2 =$$

$$= \mathbb{E} \left\| \frac{1}{m} \sum_{i=1}^m \underbrace{g(x_u, \xi_{k,i}) - f'(x_u)}_{\delta_i} \right\|^2 \quad \text{①}$$

$\mathbb{E} \delta_i = 0$

$$\text{①} \quad \mathbb{E} \left\langle \frac{1}{m} \sum_{i=1}^m \delta_i, \frac{1}{m} \sum_{j=1}^m \delta_j \right\rangle = \frac{1}{m^2} \sum_{i=1}^m \mathbb{E} \|\delta_i\|^2 + \frac{2}{m^2} \sum_{i < j} \overbrace{\mathbb{E} \langle \delta_i, \delta_j \rangle}^{=0}$$

$$\leq \frac{1}{m} \sigma^2.$$

$$\frac{LR_0^2}{\epsilon} \stackrel{?}{=} \left[ \frac{\sigma R_0}{\epsilon \sqrt{m}} \right]^2 = \frac{\sigma^2 R_0^2}{m \epsilon^2}$$

$$\Rightarrow m = 1 + \frac{\sigma^2 R_0^2}{\epsilon^2} \cdot \frac{\epsilon}{LR_0^2} = \frac{\sigma^2}{\epsilon L} + 1.$$

The iteration complexity:  $\leftarrow$  solving some subproblem

$$O\left(\frac{LR_0^2}{\epsilon}\right)$$

The total # samples:

$$O\left(\frac{LR_0^2}{\epsilon}\right) * m = O\left(\frac{LR_0^2}{\epsilon} + \frac{\sigma^2 R_0^2}{\epsilon^2}\right)$$

$$m \approx 100$$

Acceleration:  $\sqrt{\frac{LR_0^2}{\epsilon}}$

## 2. SURG

Finite-Sum Minimization:

$$f(x) := \frac{1}{N} \sum_{i=1}^N f_i(x), \quad f_i(x) \text{ - Lipschitz, convex}$$

We use in the algorithm: sample  $i$

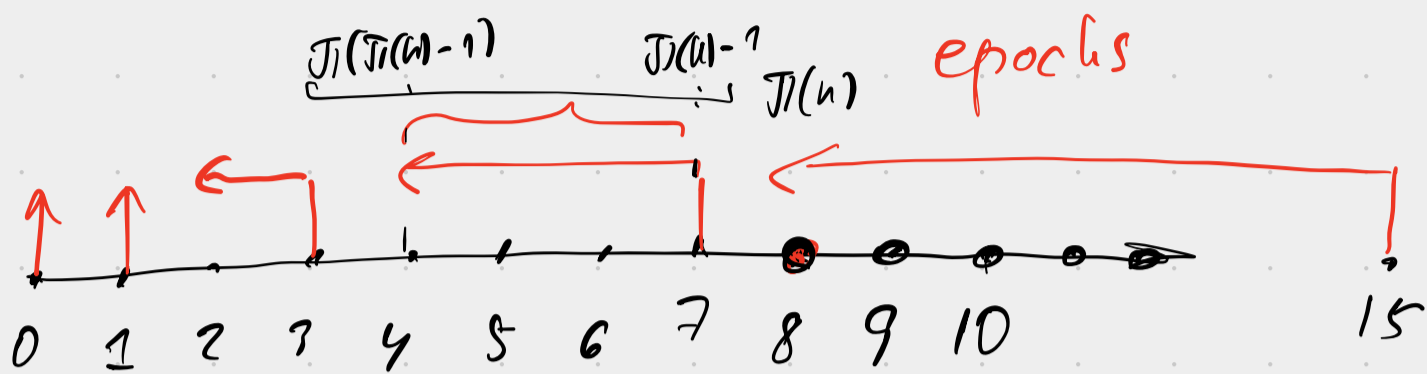
$$g_u = f_i'(x_u) - f_i'(z_u) + \underline{f'(z_u)}$$

$z_u$  - fixed point from the past.  $\leftarrow$  compute it

Let  $\pi(k)$  the largest power of 2 that is less than or equal to  $k$ .

$$\pi(k) = 2^{\lfloor \log_2 k \rfloor}$$

$$\pi(0) = 0$$



$$z_k = X_{\pi(k)} \quad ?$$

$$z_k := \frac{1}{\pi(k) - \pi(\pi(k) - 1)} \sum_{i=\pi(\pi(k) - 1)}^{\pi(k) - 1} X_i$$

lemma  $\mathbb{E} g_n = f'(x_n)$

$$\mathbb{E} \|g_n - f'(x_n)\|^2 \leq 6L \mathbb{E}[f(x_n) - f^*] + 6L \mathbb{E}[f(z_n) - f^*]$$

$\rightarrow 0$