

From the previous lecture:

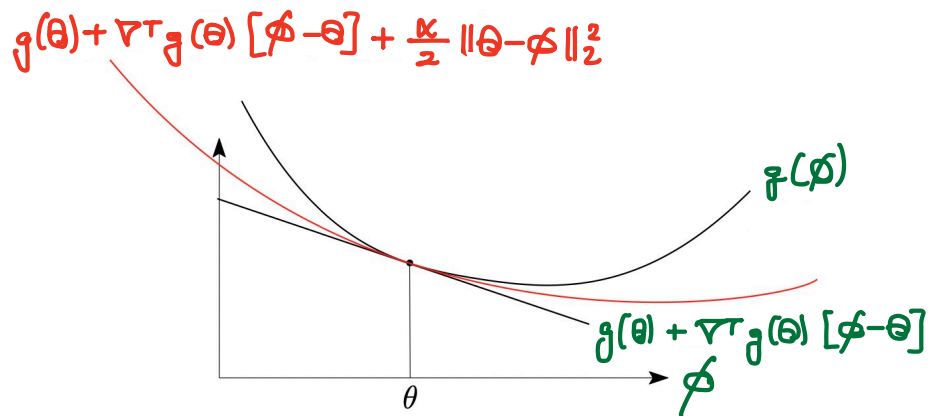
If  $g: \mathbb{R}^d \rightarrow \mathbb{R}$  is convex and Lipschitz cont., then **projected subgradient descent** achieves  $g\left(\frac{1}{T} \sum_{t=0}^{T-1} \theta_t\right) - g(\theta^*) \leq \frac{1}{\sqrt{T}}$ , which implies  $\mathcal{O}\left(\frac{1}{\varepsilon^2}\right)$  oracle complexity.

**Today**: improved convergence rates under strong convexity and smoothness.

DEF 1 ( $\alpha$ -strong convexity)

$g: \mathbb{R}^d \rightarrow \mathbb{R}$  differentiable on  $\text{dom}(g)$  is  $\alpha$ -strongly convex ( $\alpha \geq 0$ ) if:

$$g(\theta) + \nabla^T g(\theta) [\theta' - \theta] \leq g(\theta') - \frac{\alpha}{2} \cdot \|\theta - \theta'\|_2^2, \quad \forall \theta, \theta' \in \text{dom}(g).$$



PROP 1 (a)  $g: \mathbb{R}^d \rightarrow \mathbb{R}$  diff. is  $\alpha$ -strongly convex iff

$$\|\nabla g(\theta) - \nabla g(\theta')\|_2 \geq \alpha \cdot \|\theta - \theta'\|_2, \quad \forall \theta, \theta' \in \text{dom} g.$$

(b)  $g: \mathbb{R}^d \rightarrow \mathbb{R}$  twice diff. is  $\alpha$ -strongly convex iff

$$\nabla^2 g(\theta) \geq \alpha I, \quad \forall \theta \in \text{dom}(g). \quad (\text{i.e., } \nabla^2 g(\theta) - \alpha I \text{ is p.s.d.})$$

Proof: Exercise.

LEMMA 1 |  $g$  is  $\alpha$ -sc. iff  $\theta \mapsto g(\theta) - \frac{\alpha}{2} \|\theta\|_2^2$  is convex.

LEMMA 2 |  $g: \mathbb{R}^d \rightarrow \mathbb{R}$  is strongly convex and continuous.  
 $\Rightarrow$  there exists a unique  $\theta^* \in \underset{\theta \in \text{dom } g}{\text{argmin}} g(\theta)$ .

The proofs of lemma 1 and lemma 2 can be found in the appendix.

The following ineq. is of fundamental importance.

LEMMA 3 (Lojasiewicz ineq.)  $g: \mathbb{R}^d \rightarrow \mathbb{R}$  differentiable,  $\alpha$ -sc.

Then,  $\|\nabla g(\theta)\|_2^2 \geq 2\alpha \cdot [g(\theta) - \min_{\theta'} g(\theta')]$ ,  $\forall \theta \in \text{dom}(g)$

**Pf:** Let  $h(\phi) = g(\theta) + \nabla^T g(\theta) [\phi - \theta] + \frac{\alpha}{2} \|\theta - \phi\|_2^2$ .

Then,  $h$  is minimized with  $\phi^* = \theta - \frac{1}{\alpha} \nabla g(\theta)$  by the 1st-order condition for optimality. Then,

$$\begin{aligned} g(\phi) &\geq g(\theta) + \nabla^T g(\theta) [\phi - \theta] + \frac{\alpha}{2} \|\phi - \theta\|_2^2 \\ &\geq g(\theta) - \frac{1}{\alpha} \|\nabla g(\theta)\|_2^2 + \frac{1}{2\alpha} \|\nabla g(\theta)\|_2^2 = g(\theta) - \frac{1}{2\alpha} \|\nabla g(\theta)\|_2^2 \end{aligned}$$

Thus,

$$\inf_{\phi} g(\phi) \geq g(\theta) - \frac{1}{2\alpha} \|\nabla g(\theta)\|_2^2.$$

$$\Rightarrow \|\nabla g(\theta)\|_2^2 \geq 2\alpha \cdot [g(\theta) - \inf_{\phi} g(\phi)], \forall \theta \in \text{dom } g$$

An immediate implication is that whenever  $\nabla g(\theta) = 0$ ,

$$g(\theta) = \inf_{\phi} g(\phi).$$

## PROPOSITION 2 (regularization)

Let  $g: \mathbb{R}^d \rightarrow \mathbb{R}$  be a convex function, and  $h: \mathbb{R}^d \rightarrow \mathbb{R}$  be an  $l$ -sc function. Then, for any  $\lambda \in \mathbb{R}_+$ ,

$$\theta \mapsto g(\theta) + \lambda h(\theta) \text{ is } \lambda\text{-sc.}$$

**Pf:**  $\varphi(\theta) = g(\theta) + \lambda h(\theta) \Rightarrow \nabla \varphi(\theta) = \nabla g(\theta) + \lambda \nabla h(\theta)$

Then,

$$\begin{aligned} \varphi(\theta) + \nabla^T \varphi(\theta) [\theta' - \theta] &= g(\theta) + \nabla^T g(\theta) [\theta' - \theta] + \lambda [h(\theta) + \nabla^T h(\theta) (\theta' - \theta)] \\ &\leq g(\theta') + \lambda \cdot [h(\theta') - \frac{1}{2} \|\theta - \theta'\|_2^2] \\ &= \varphi(\theta') - \frac{\lambda}{2} \|\theta - \theta'\|_2^2, \quad \forall \theta, \theta' \in \text{dom } g \cap \text{dom } h \end{aligned}$$

Note: A canonical example is Tikhonov regularization:

$$g(\theta) + \frac{\lambda}{2} \|\theta\|_2^2, \quad \lambda > 0 \text{ is } \lambda\text{-sc.}$$

## GRADIENT DESCENT FOR STRONGLY CONVEX FUNCTIONS

### Algorithm: Gradient Descent

Inputs:  $(\gamma_t)_{t \geq 0}$  step-sizes

$\theta_0 \in \text{dom } g$

for  $t=0, 1, \dots, T-1$

$$\theta_{t+1} = \theta_t - \gamma_t \nabla g(\theta_t)$$

what does GD do? Proximal form:

$$\theta_{t+1} \in \arg \min_{\theta} \left\{ \underbrace{g(\theta_t) + \nabla^T g(\theta_t) [\theta - \theta_t]}_{\text{linear approximation around } g(\theta_t)} + \frac{1}{2\gamma} \|\theta - \theta_t\|_2^2 \right\}$$

linear approximation  
around  $g(\theta_t)$

$\|\theta - \theta_t\|_2^2$  penalizes moving from  $\theta_t$ .

large step-size  $\Rightarrow$  small penalty

$\Rightarrow$  move far from  $\theta_t$ .

More generally,

$$\theta_{t+1} \in \arg \min_{\theta} \left\{ g(\theta_t) + \nabla^T g(\theta_t) [\theta - \theta_t] + \frac{1}{\gamma} \underline{D(\theta, \theta_t)} \right\}$$

Bregman divergence

Mirror descent for non-Euclidean opt. (Nemirovski, Yudin, 1981)

# THEOREM 1 (Convergence of GD for $\alpha$ -sc., $L$ -Lipschitz functions)

$g: \mathbb{R}^d \rightarrow \mathbb{R}$  differentiable,  $\alpha$ -sc. and  $L$ -Lipschitz.

Then, for any  $T \geq 1$ , with  $\eta_t = \frac{1}{\alpha(t+1)}$ ,

$$(a) \quad \|\theta_T - \theta^*\|_2^2 \leq \frac{L^2 (1 + \log T)}{\beta^2 T},$$

$$(b) \quad g\left(\frac{1}{T} \sum_{t=0}^{T-1} \theta_t\right) - g(\theta^*) \leq \frac{L^2}{2\beta T} [1 + \log T],$$

where  $\theta^*$  is the unique minimizer of  $g$ .

**Pf:** Lyapunov function :  $L(\theta) = \|\theta - \theta^*\|_2^2$ ,  $\theta \in \mathbb{R}^d$ .

Then, by  $\alpha$ -sc. and  $L$ -Lipschitz continuity,  $\forall t \geq 0$ ,

$$\begin{aligned} L(\theta_{t+1}) &= L(\theta_t) - 2\eta_t \nabla^T g(\theta_t) [\theta_t - \theta^*] + \eta_t^2 \|\nabla g(\theta_t)\|_2^2 \\ &\leq L(\theta_t) - 2\eta_t \left[ g(\theta_t) - g(\theta^*) + \frac{\alpha}{2} \|\theta_t - \theta^*\|_2^2 \right] + \eta_t^2 L^2 \\ &= L(\theta_t) [1 - \eta_t \alpha] - 2\eta_t \Delta_t + \eta_t^2 L^2, \quad \Delta_t := g(\theta_t) - g(\theta^*). \end{aligned}$$

(a) Since  $\Delta_t \geq 0$ ,  $\forall t$ , using  $\eta_t = \frac{1}{\alpha(t+1)}$ ,

$$\begin{aligned} L(\theta_T) &\leq \frac{T-1}{T} L(\theta_{T-1}) + \frac{L^2}{\alpha^2 T^2} \leq \frac{T-1}{T} \left[ \frac{T-2}{T-1} L(\theta_{T-2}) + \frac{L^2}{\alpha^2 (T-1)^2} \right] + \frac{L^2}{\alpha^2 T^2} \\ &\leq \frac{T-2}{T-1} L(\theta_{T-2}) + \frac{L^2}{\alpha^2 T} \left[ \frac{1}{T} + \frac{1}{T-1} \right] \\ &\vdots \\ &\leq \frac{L^2}{\alpha^2 T} \left[ 1 + \frac{1}{2} + \dots + \frac{1}{T} \right] \leq \frac{L^2}{\alpha^2 T} [\log T + 1]. \end{aligned}$$

$$(b) \quad L(\theta_T) \leq \frac{T-1}{T} L(\theta_{T-1}) - \frac{2}{\alpha T} \Delta_{T-1} + \frac{L^2}{\alpha^2 T^2}$$

$$\leq -\frac{2}{\alpha T} \sum_{t=0}^{T-1} \Delta_t + \frac{L^2}{\alpha^2 T} (1 + \log T) \quad \text{by induction.}$$

$$\text{Then, } \frac{1}{T} \sum_{t=0}^{T-1} g(\theta_t) - g(\theta^*) \leq \frac{L^2}{2\alpha T} [1 + \log T] \quad \text{since } L \geq 0.$$

$$\text{Jensen's ineq. } \Rightarrow g\left(\frac{1}{T} \sum_{t=0}^{T-1} \theta_t\right) - g(\theta^*) \leq \frac{1}{T} \sum_{t=0}^{T-1} g(\theta_t) - g(\theta^*) \leq \frac{L^2 (1 + \log T)}{2\alpha T}$$

Remarks: (1) No need for projection when we have sc.

(2) Average-iterate rate:  $g(\bar{\theta}_T) - g^* = \tilde{O}\left(\frac{1}{T}\right)$ ,  
without sc.  $g(\bar{\theta}_T) - g^* = \tilde{O}\left(\frac{1}{\sqrt{T}}\right)$ .

(3) Last-iterate rate for parameter convergence:  
 $\|\theta_T - \theta^*\|_2^2 = \tilde{O}\left(\frac{1}{T}\right)$ .

PP (LEMMA 1): Let  $\varphi(\theta) = f(\theta) - \frac{\alpha}{2} \|\theta\|_2^2$ ,  $\theta \in \text{dom } f$ . Then,

$$\nabla \varphi(\theta) = \nabla f(\theta) - \alpha \theta.$$

$$\begin{aligned} (\Rightarrow) \quad \varphi(\theta) + \nabla^T \varphi(\theta) [\theta' - \theta] &= f(\theta) - \frac{\alpha}{2} \|\theta\|_2^2 + [\nabla f(\theta) - \alpha \theta]^T (\theta' - \theta) \\ &= f(\theta) + \nabla^T f(\theta) [\theta' - \theta] + \frac{\alpha}{2} \|\theta\|_2^2 - \alpha \theta^T \theta' \\ &\leq f(\theta') + \frac{\alpha}{2} \|\theta\|_2^2 - \alpha \theta^T \theta' - \frac{\alpha}{2} \|\theta - \theta'\|_2^2 \\ &= f(\theta') - \frac{\alpha}{2} \|\theta'\|_2^2 = \varphi(\theta'). \end{aligned}$$

( $f$  is  $\alpha$ -sc.)

$$(\Leftarrow) \quad f(\theta) - \frac{\alpha}{2} \|\theta\|_2^2 + (\nabla f(\theta) - \alpha \theta)^T (\theta' - \theta) \leq f(\theta') - \frac{\alpha}{2} \|\theta'\|_2^2$$

Then,

$$\begin{aligned} f(\theta) + \nabla^T f(\theta) [\theta' - \theta] &\leq f(\theta') - \frac{\alpha}{2} \|\theta'\|_2^2 + \alpha \theta^T (\theta' - \theta) + \frac{\alpha}{2} \|\theta\|_2^2 \\ &= f(\theta') - \frac{\alpha}{2} \|\theta'\|_2^2 + \alpha \theta^T \theta' - \frac{\alpha}{2} \|\theta\|_2^2 \\ &= f(\theta') - \frac{\alpha}{2} \|\theta' - \theta\|_2^2. \quad \blacksquare \end{aligned}$$