

Recall: $\mathcal{H} \subset \mathbb{R}^d$, $f: \mathcal{H} \rightarrow \mathbb{R}$. Suppose that f is differentiable over its domain. Then,

∇f if f is convex, then

$$f(\theta) + \nabla f(\theta)^T [\theta' - \theta] \leq f(\theta'), \quad \theta, \theta' \in \mathcal{H}$$

Non-Differentiable Convex Optimization

Note that a convex function is not always differentiable (everywhere) on its domain. E.g., $f(\theta) = |\theta|$, $f(\theta) = \max\{0, \theta\}$.
ReLU

Next, we generalize gradients.

DEF (Subgradients) $u \in \mathbb{R}^d$ is a subgradient of f at $\theta \in \mathcal{H}$

if $f(\theta') \geq f(\theta) + u^T [\theta' - \theta]$ for any $\theta' \in \mathcal{H}$.

The set of subgradients of f at $\theta \in \mathcal{H}$ is called the subdifferential and is denoted as $\partial f(\theta)$.

PROP 1 (Existence of subgradients)

Let $\mathcal{H} \subset \mathbb{R}^d$ be convex, and $f: \mathcal{H} \rightarrow \mathbb{R}$.

(a) If $\partial f(\theta) \neq \emptyset$ for any $\theta \in \mathcal{H}$, then f is convex.

(b) If f is convex, then $\partial f(\theta) \neq \emptyset$ for any $\theta \in \mathcal{H}$.

Proof of PROP 1:

(a) Let $u \in \partial f(\gamma\theta + (1-\gamma)\theta')$. Then, by definition,

$$\begin{aligned} 1-\gamma / & \quad f(\theta') \geq f(\gamma\theta + (1-\gamma)\theta') + \gamma \cdot u^T (\theta' - \theta) \\ \gamma / & \quad f(\theta) \geq f(\gamma\theta + (1-\gamma)\theta') - (1-\gamma) u^T (\theta' - \theta) \end{aligned}$$

$$\gamma f(\theta) + (1-\gamma) f(\theta') \geq f(\gamma\theta + (1-\gamma)\theta'). \Rightarrow f \text{ is convex.}$$

(b) We will construct a subgradient of f at $\theta \in \text{int dom } f$. For this, we will use the supporting hyperplane theorem (see appendix).

Let $\theta \in \text{int}(\mathcal{H})$. Then, $(\theta, f(\theta)) \in \text{epi}(f)$, indeed $(\theta, f(\theta)) \in \text{bd epi}(f)$. Since f is convex, $\text{epi}(f)$ is convex. Thus, by the supporting hyperplane theorem, $\exists (w, v) \in \mathbb{R}^d \times \mathbb{R}$ s.t. $(w, v) \neq 0$, and

$$[w^T \ v] \begin{bmatrix} \theta \\ f(\theta) \end{bmatrix} \geq [w^T \ v] \begin{bmatrix} \theta' \\ t \end{bmatrix} \text{ for any } (\theta', t) \in \text{epi}(f).$$

By tending $t \rightarrow \infty$, we see that $v \leq 0$ must hold.

Since $\theta \in \text{int}(\mathcal{H})$, for sufficiently small $\varepsilon > 0$,

$$\tilde{\theta} = \theta + \varepsilon \omega \in \mathcal{H}.$$

$\Rightarrow v \neq 0$, since otherwise

$$\begin{bmatrix} \omega^T & 0 \end{bmatrix} \begin{bmatrix} \theta \\ g(\theta) \end{bmatrix} \geq \begin{bmatrix} \omega^T & 0 \end{bmatrix} \begin{bmatrix} \theta + \varepsilon \omega \\ t \end{bmatrix}$$

$$\Rightarrow \omega^T \theta \geq \omega^T \theta + \varepsilon \|\omega\|^2 \Rightarrow \omega = 0, \text{ which would contradict with } (\omega, v) \neq 0.$$

Let $t = g(\theta')$. Then, for some $v < 0$,

$$\begin{bmatrix} \omega^T & v \end{bmatrix} \begin{bmatrix} \theta \\ g(\theta) \end{bmatrix} \geq \begin{bmatrix} \omega^T & v \end{bmatrix} \begin{bmatrix} \theta' \\ g(\theta') \end{bmatrix}$$

$$\Rightarrow \omega^T(\theta - \theta') + v g(\theta) \geq v g(\theta'), \quad \theta' \in \mathcal{H}.$$

$$\Rightarrow g(\theta) + \frac{\omega^T}{v}(\theta - \theta') \leq g(\theta'), \quad \text{since } v < 0.$$

Thus, $\partial g(\theta) \neq \emptyset$ for $\theta \in \text{int}(\mathcal{H})$. ■

Remark: If g is differentiable, (b) \neq PROP 1 automatically holds since $\nabla g(\theta) \in \partial g(\theta)$ by the first-order condition for convexity.

Pf | (Existence of subgradients : detailed)

(A) $\partial f(\theta) \neq \emptyset, \forall \theta \Rightarrow f$ is convex.

(B) f is convex $\Rightarrow \partial f(\theta) \neq \emptyset - \forall \theta \in \mathbb{H}$.

Pf (B) Supporting hyperplane theorem:

Let $\mathbb{H} \subset \mathbb{R}^d$ be a convex set, and $\theta_0 \in \text{bd } \mathbb{H}$.
Then, $\exists w \in \mathbb{R}^d$ s.t. $w \neq 0$ - and
 $w^T \theta \leq w^T \theta_0, \forall \theta \in \mathbb{H}$.

f is convex $\Rightarrow \text{epi}(f)$ is convex.

$\theta \in \text{mt}(\mathbb{H})$. Then, $(\theta, f(\theta)) \in \text{epi}(f)$, $(\theta, f(\theta)) \in \text{bd epi}(f)$

By SHT, $\exists (w, v) \in \mathbb{R}^d \times \mathbb{R}$ s.t. $(w, v) \neq 0$ and

$$[w^T \ v] \begin{bmatrix} \theta \\ f(\theta) \end{bmatrix} \geq [w^T \ v] \begin{bmatrix} \theta' \\ t \end{bmatrix}$$

for any $\begin{bmatrix} \theta' \\ t \end{bmatrix} \in \text{epi}(f)$. Then, since $t \geq f(\theta')$, we

conclude that $v \leq 0$. Since $\theta \in \text{mt}(\mathbb{H})$,

$$\exists \varepsilon > 0 \text{ s.t. } \theta + \varepsilon w \in \mathbb{H}.$$

Thus, if $v = 0$,

$$[w^T \ 0] \begin{bmatrix} \theta \\ f(\theta) \end{bmatrix} \geq [w^T \ 0] \begin{bmatrix} \theta + \varepsilon w \\ t \end{bmatrix}$$

$$\Rightarrow w = 0 \Rightarrow \Leftarrow$$

Hence, $v < 0$ must hold.

Then, let $t = f(\theta')$.

$$\Rightarrow w^T \theta + v f(\theta) \geq w^T \theta' + v f(\theta')$$

$$\Rightarrow f(\theta) + \frac{w^T (\theta - \theta')}{v} \leq f(\theta')$$

$$f(\theta) + u^T (\theta' - \theta) \leq f(\theta')$$

$$\text{with } u = -\frac{w}{v}.$$

DEF | (Lipschitz continuity) Let $\mathcal{H} \subset \mathbb{R}^d$.

$g: \mathcal{H} \rightarrow \mathbb{R}$ is Lipschitz continuous if there exists a positive constant L s.t.

$$|g(\theta) - g(\theta')| \leq L \cdot \|\theta - \theta'\|_2, \quad \forall \theta, \theta' \in \mathcal{H}.$$

Exercise: If $g: \mathcal{H} \rightarrow \mathbb{R}^d$ is differentiable and $\|\nabla g(\theta)\|_2 \leq L$ for any $\theta \in \mathcal{H}$, then g is L -Lipschitz continuous.

THEOREM 1 | (Lipschitz continuity of convex functions)

Let $g: \mathbb{R}^d \rightarrow \mathbb{R}$, and K be a compact set contained in $\text{int dom}(g)$. Then, g is Lipschitz continuous on K .

For the proof of Theorem 1, we will need the following lemmas.

Lemma 1 | (Convex functions are locally bounded)

Let g be convex, and $\theta_0 \in \text{int dom}(g)$. Then, g is locally bounded:
 $\exists \varepsilon > 0$ and $M(\theta_0, \varepsilon) > 0$ s.t.

$$\max_{\theta \in B_2(\theta_0, \varepsilon)} g(\theta) \leq M(\theta_0, \varepsilon).$$

Lemma 2 | (Convex functions are locally Lipschitz)

Let g be convex, and $\theta_0 \in \text{int dom}(g)$. Then, g is locally Lipschitz:
 $\exists \varepsilon > 0$ and $L(\theta_0, \varepsilon) < \infty$ s.t.

$$|g(\theta) - g(\theta_0)| \leq L(\theta_0, \varepsilon) \cdot \|\theta - \theta_0\|_2, \quad \forall \theta \in B_2(\theta_0, \varepsilon).$$

Lemma 3 | The following statements hold for a convex function:

(a) If $\theta_0 \in \text{int dom}(g)$, and g is locally Lipschitz, then
 $\|u\| \leq L(\theta_0, \varepsilon) < \infty$ for any $u \in \partial g(\theta_0)$

(b) If $\|u\|_2 \leq L(\theta_0, \varepsilon)$ for $u \in \partial g(\theta_0)$, then
 $g(\theta_0) - g(\theta) \leq L(\theta_0, \varepsilon) \cdot \|\theta_0 - \theta\|_2.$

PROP 2 | (Local minima are global minima under convexity)

Let g be convex.

(i) If θ^* is a local minimum of g , then θ^* is a global minimum of g .

(ii) The above happens if and only if $0 \in \partial g(\theta^*)$.

Proof: Clearly, $0 \in \partial g(\theta) \iff \theta$ is a global minimum of g .

Now, assume that θ is a local minimum of g :

$$\exists \varepsilon > 0 \text{ s.t. } g(\theta') \geq g(\theta), \forall \theta' \in B_2(0, \varepsilon).$$

Then, for small enough γ , for any $\tilde{\theta} \in \text{dom } g$,

$$g(\theta) \leq g((1-\gamma)\theta + \gamma\tilde{\theta}) \leq (1-\gamma)g(\theta) + \gamma g(\tilde{\theta})$$

$$\Rightarrow \gamma g(\theta) \leq \gamma g(\tilde{\theta}) \Rightarrow g(\theta) \leq g(\tilde{\theta}). \quad \blacksquare$$

PROP 3 | (First-order optimality condition : constrained case)

Let g be convex, and $\mathcal{H} \subset \mathbb{R}^d$ a closed convex set on which g is differentiable. Then,

$$\theta^* \in \underset{\theta \in \mathcal{H}}{\text{argmin}} g(\theta),$$

if and only if

$$\nabla^T g(\theta^*) [\theta^* - \theta] \leq 0, \quad \forall \theta \in \mathcal{H}.$$

Proof: (\Leftarrow) $g(\theta) + \nabla^T g(\theta^*) [\theta - \theta^*] \leq g(\theta)$, $\forall \theta \in \mathcal{H}$ by convexity.

$$g(\theta) - g(\theta^*) \geq -\nabla^T g(\theta^*) [\theta^* - \theta] \geq 0, \quad \forall \theta \in \mathcal{H}.$$

$$\Rightarrow \theta^* \in \underset{\theta \in \mathcal{H}}{\text{argmin}} g(\theta).$$

(\Rightarrow) For $\theta \in \mathcal{H}$, let $h(t) := g(\theta^* + t[\theta - \theta^*])$, $t \in [0, 1]$.

h is a convex function. Pf: Exercise.

Then,
 $\frac{dh(t)}{dt} = \nabla^T g(\theta^* + t[\theta - \theta^*]) \cdot [\theta - \theta^*]$. Suppose to the contrary that $\frac{dh(t)}{dt}|_{t=0} < 0$. Then, $\exists t_0 \in (0, 1)$ small enough s.t. $h(t_0) < h(0)$. This is a contradiction since $\theta^* \in \underset{\theta \in \mathcal{H}}{\text{argmin}} g(\theta)$.

NONDIFFERENTIABLE CONVEX OPTIMIZATION

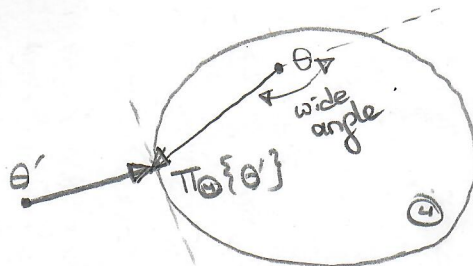
DEF (Projection) Let $\mathcal{H} \subset \mathbb{R}^d$ be a compact and convex set. For any $\theta \in \mathbb{R}^d$, let

$$\Pi_{\mathcal{H}}\{\theta\} = \operatorname{argmin}_{\theta' \in \mathcal{H}} \|\theta - \theta'\|_2.$$

LEMMA | Let $\theta \in \mathcal{H}$, and $\theta' \in \mathbb{R}^d$. Then,

$$\left(\Pi_{\mathcal{H}}\{\theta'\} - \theta \right)^T \left(\Pi_{\mathcal{H}}\{\theta'\} - \theta' \right) \leq 0.$$

Remark (Geometric intuition)



Proof: Let $h(\tilde{\theta}) = \frac{1}{2} \|\tilde{\theta} - \theta'\|_2^2$, $\tilde{\theta} \in \mathbb{R}^d$. Then, obviously, h is a convex function. By the first-order condition for optimality (PROP 3),

$$\nabla^T h(\Pi_{\mathcal{H}}\{\theta'\}) [\Pi_{\mathcal{H}}\{\theta'\} - \tilde{\theta}] \leq 0, \quad \forall \tilde{\theta} \in \mathcal{H}.$$

Thus,

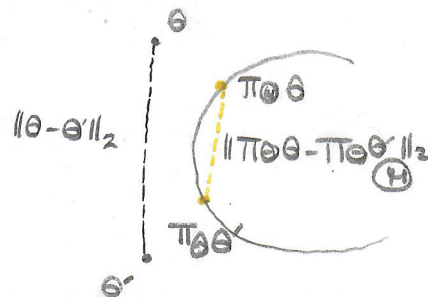
$$[\Pi_{\mathcal{H}}\{\theta'\} - \theta']^T [\Pi_{\mathcal{H}}\{\theta'\} - \tilde{\theta}] \leq 0, \quad \forall \tilde{\theta} \in \mathcal{H}.$$

PROP 4 | (Projection is a nonexpansive operator)

Let $\mathcal{H} \subset \mathbb{R}^d$ be a compact and convex set. Then, for any $\theta, \theta' \in \mathbb{R}^d$,

$$\|\Pi_{\mathcal{H}}\{\theta\} - \Pi_{\mathcal{H}}\{\theta'\}\|_2 \leq \|\theta - \theta'\|_2.$$

Proof: By LEMMA above,



Proof:
$$\begin{aligned} & [\pi_{\mathcal{H}}(\theta) - \theta]^T [\pi_{\mathcal{H}}(\theta) - \underbrace{\pi_{\mathcal{H}}(\theta')}_{\in \mathcal{H}}] \leq 0 \\ & [\pi_{\mathcal{H}}(\theta') - \theta']^T [\pi_{\mathcal{H}}(\theta') - \pi_{\mathcal{H}}(\theta)] \leq 0 \end{aligned}$$

Thus,

$$\begin{aligned} & [\pi_{\mathcal{H}}(\theta) - \pi_{\mathcal{H}}(\theta')]^T [\pi_{\mathcal{H}}(\theta) - \theta - (\pi_{\mathcal{H}}(\theta') - \theta')] \\ &= \|\pi_{\mathcal{H}}(\theta) - \pi_{\mathcal{H}}(\theta')\|_2^2 - (\pi_{\mathcal{H}}(\theta) - \pi_{\mathcal{H}}(\theta'))^T (\theta' - \theta) \leq 0 \end{aligned}$$

Hence, by using Cauchy-Schwarz,

$$\|\pi_{\mathcal{H}}(\theta) - \pi_{\mathcal{H}}(\theta')\|_2 \leq \|\theta - \theta'\|_2.$$

ALGORITHM 1. $\mathcal{H} \subset \mathbb{R}^d$ compact and convex set

Inputs : $\theta_0 \in \mathcal{H}$, $\gamma > 0$,

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

$$\tilde{\theta}_{t+1} = \theta_t - \gamma u_t, \quad u_t \in \partial g(\theta_t),$$

$$\theta_{t+1} = \pi_{\mathcal{H}}(\tilde{\theta}_{t+1}).$$

Return :
$$\bar{\theta}_T = \frac{1}{T} \sum_{t=0}^{T-1} \theta_t.$$

THEOREM (Lyapunov - deterministic)

Let $\mathcal{H} \subset \mathbb{R}^d$ be a compact and convex set.

For a given sequence $\theta_t \in \mathcal{H}$, $t \in \mathbb{Z}_+$, assume that there exist

$$L : \mathcal{H} \rightarrow \mathbb{R}_+,$$

$$c : \mathcal{H} \rightarrow \mathbb{R}_+ \text{ (convex),}$$

$$\varepsilon \geq 0 \text{ and } V \in \mathbb{R} \text{ s.t.}$$

for any $\eta > 0$, $t \geq 0$,

$$L(\theta_{t+1}) - L(\theta_t) \leq -2\eta c(\theta_t) + \eta \cdot \varepsilon + \eta^2 V^2$$

Then, for any $T \geq 1$:

$$c\left(\frac{1}{T} \sum_{t=0}^{T-1} \theta_t\right) \leq \frac{L(\theta_0)}{2\eta T} + \frac{\varepsilon}{2} + \frac{\eta V^2}{2}$$

COROLLARY With $\eta = \sqrt{\frac{L(\theta_0)}{T}} \cdot \frac{1}{V}$, we have:

$$c\left(\frac{1}{T} \sum_{t=0}^{T-1} \theta_t\right) \leq \sqrt{\frac{L(\theta_0)}{T}} + \frac{\varepsilon}{2}.$$

Proof By telescoping sum over $t=0, 1, \dots, T-1$:

$$L(\theta_T) - L(\theta_0) \leq -2\eta \sum_{t=0}^{T-1} c(\theta_t) + \eta \cdot T \cdot \varepsilon + \eta^2 \cdot T \cdot V^2$$

Rearranging terms, since $L(\theta_T) \geq 0$,

$$\frac{1}{T} \sum_{t=0}^{T-1} c(\theta_t) \leq \frac{L(\theta_0)}{2\eta T} + \frac{\varepsilon}{2} + \frac{\eta V^2}{2}.$$

Since c is a convex function, by Jensen's inequality,

$$c\left(\frac{1}{T} \sum_{t=0}^{T-1} \theta_t\right) \leq \frac{1}{T} \sum_{t=0}^{T-1} c(\theta_t).$$

From the above inequalities, the result follows.

Performance analysis of Projected Subgradient Descent:

THEOREM (Near-optimality of Proj-GD) suppose $\sup_{\theta, \theta' \in \mathcal{H}} \|\theta - \theta'\|_2 \leq \rho$.

Algorithm Proj-GD with $\gamma = \frac{\rho}{L\sqrt{T}}$ satisfies

$$f\left(\frac{1}{T} \sum_{t=0}^{T-1} \theta_t\right) - \min_{\theta \in \mathcal{H}} f(\theta) \leq \frac{\rho L}{\sqrt{T}}$$

for any $T \geq 1$.

Pf: using $\mathcal{L}(\theta) = \|\theta - \theta^*\|_2^2$ where $\theta^* \in \arg\min_{\theta \in \mathcal{H}} f(\theta)$:

$$\begin{aligned} \mathcal{L}(\theta_{t+1}) &= \|\theta_{t+1} - \theta^*\|_2^2 \\ &= \|\Pi_{\mathcal{H}}(\tilde{\theta}_{t+1}) - \Pi_{\mathcal{H}}(\theta^*)\|_2^2 \quad \text{since } \theta^* \in \mathcal{H} \\ &\leq \|\tilde{\theta}_{t+1} - \theta^*\|_2^2 \quad \text{by nonexpansiveness of } \Pi_{\mathcal{H}} \text{ (PROP 4)} \\ &= \|\theta_t - \gamma u_t - \theta^*\|_2^2 \\ &= \|\theta_t - \theta^*\|_2^2 - 2\gamma u_t^T(\theta_t - \theta^*) + \gamma^2 \|u_t\|_2^2 \\ &\leq \mathcal{L}(\theta_t) - 2\gamma u_t^T(\theta_t - \theta^*) + \gamma^2 L^2 \end{aligned}$$

Then,

$$\mathcal{L}(\theta_{t+1}) - \mathcal{L}(\theta_t) \leq -2\gamma \cdot u_t^T(\theta_t - \theta^*) + \gamma^2 L^2.$$

By the defn of subgradient,

$$f(\theta_t) + u_t^T(\theta^* - \theta_t) \leq f(\theta^*).$$

Hence,

$$u_t^T(\theta_t - \theta^*) - f(\theta_t) \geq -f(\theta^*)$$

$$\Rightarrow u_t^T(\theta_t - \theta^*) \geq c(\theta_t)$$

$$\text{where } c(\theta_t) = f(\theta_t) - f(\theta^*).$$

By Theorem (Lyapunov), since $\mathcal{L} \geq 0$ and $\theta \mapsto c(\theta)$ is convex, we have:

$$c\left(\frac{1}{T} \sum_{t=0}^{T-1} \theta_t\right) \leq \frac{L \sqrt{\|\theta_0 - \theta^*\|_2^2}}{\sqrt{T}} \quad \text{since } \max_{\theta, \theta' \in \mathcal{H}} \|\theta - \theta'\|_2 \leq \rho,$$

the result holds.