

CONVEX OPTIMIZATION

Why? Although almost all deep learning problems are highly nonconvex, most of the tools developed for convex optimization are directly or indirectly used in deep learning.

DEF | (Convex set, convex function)

$\Theta \subset \mathbb{R}^d$ is a convex set if $\theta, \theta' \in \Theta \Rightarrow \gamma\theta + (1-\gamma)\theta' \in \Theta, \forall \gamma \in [0,1]$

$f: \mathbb{R}^d \rightarrow \mathbb{R}$ is a convex function if $\text{dom}(f)$ is a convex set, and
 $f(\gamma\theta + (1-\gamma)\theta') \leq \gamma f(\theta) + (1-\gamma)f(\theta'), \forall \theta, \theta' \in \text{dom}(f), \gamma \in [0,1]$.

As simple examples, $f(\theta) = a^T \theta + b$, $g(\theta) = \|\theta\|_p$ for $p \geq 1$,
 $g(\theta) = \max\{0, 1 - a^T \theta\}$.

PROP | (epigraph) Let $\text{epi}(f) = \{(\theta, t) \in \mathbb{R}^d \times \mathbb{R} : \theta \in \text{dom}(f), t \geq f(\theta)\}$.

Then, f is convex if and only if $\text{epi}(f)$ is convex.

Pf: Exercise.

Recall that $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is differentiable if $\text{dom}(f)$ is an open set, and
 $\nabla_{\theta} f(\theta) = \begin{bmatrix} \frac{\partial f}{\partial \theta_1} \\ \vdots \\ \frac{\partial f}{\partial \theta_d} \end{bmatrix}$ exist for all $\theta \in \text{dom}(f)$.

PROP | (1st-order condition for convexity)

Suppose that $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is differentiable with convex $\text{dom}(f)$.

Then,

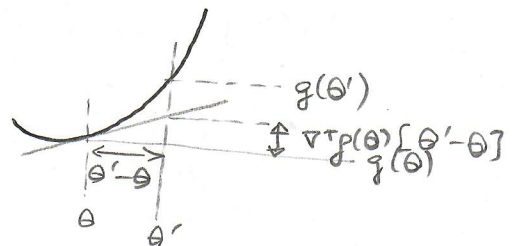
if f is convex, $f(\theta) + \nabla^T f(\theta) [\theta' - \theta] \leq f(\theta'), \forall \theta, \theta' \in \text{dom}(f)$.

Proof: For any $\theta, \theta' \in \text{dom}(f)$, $\gamma \in [0,1]$,

$$f(\gamma\theta' + (1-\gamma)\theta) \leq \gamma f(\theta') + (1-\gamma)f(\theta)$$

Rearranging terms,

$$\frac{f(\theta + \gamma[\theta' - \theta]) - f(\theta)}{\gamma} + f(\theta) \leq f(\theta').$$



f is diff. \Rightarrow as $\gamma \rightarrow 0$, $\nabla^T f(\theta) [\theta' - \theta] + f(\theta) \leq f(\theta')$.

Note: The converse of the above is also true: if $\forall \theta, \theta' \in \text{dom}(f)$, $f(\theta) + \nabla^T f(\theta) [\theta' - \theta] \leq f(\theta')$, then f is convex.

PROP | (2nd-order condition for convexity)

Suppose that $g: \mathbb{R}^d \rightarrow \mathbb{R}$ is twice-differentiable, $\text{dom } g$ is convex.

Then,

g is convex $\Leftrightarrow \nabla^2 g(\theta)$ is positive definite for all $\theta \in \text{dom } g$.

PROP | (Some convexity-preserving operations)

(1) $h: \mathbb{R}^m \rightarrow \mathbb{R}$ is convex

$\Rightarrow g(\theta) = h(A\theta + b)$ is convex for any $A \in \mathbb{R}^{m \times d}$, $b \in \mathbb{R}^m$.

(2) g_1, g_2 convex $\Rightarrow g_1 + g_2$ is convex.

g_1, g_2, \dots, g_m convex $\Rightarrow g(\theta) = \max_{i=1,2,\dots,m} g_i(\theta)$, $\theta \in \bigcap_{i=1}^m \text{dom}(g_i)$ is convex.

(3) $g: \mathbb{R}^d \rightarrow \mathbb{R}$, $h: \mathbb{R} \rightarrow \mathbb{R}$. Let $\varphi = h \circ g$.

φ is convex \Leftrightarrow $\begin{cases} g \text{ is convex, } h \text{ is convex, } h \text{ is nondecreasing} \\ g \text{ is concave, } h \text{ is convex, } h \text{ is nonincreasing} \end{cases}$

Prf: (1) $g(\gamma\theta + (1-\gamma)\theta') = h(A(\gamma\theta + (1-\gamma)\theta') + b)$
 $= h(\gamma(A\theta + b) + (1-\gamma)(A\theta' + b))$
 $\leq \gamma h(A\theta + b) + (1-\gamma)h(A\theta' + b)$ by convexity of h .
 $= \gamma g(\theta) + (1-\gamma)g(\theta')$. $\Rightarrow g$ is convex.

(2) $g_i(\gamma\theta + (1-\gamma)\theta') \leq \gamma g_i(\theta) + (1-\gamma)g_i(\theta')$, $i=1,2$.

Then, $(g_1 + g_2)(\gamma\theta + (1-\gamma)\theta') \leq \gamma(g_1 + g_2)(\theta) + (1-\gamma)(g_1 + g_2)(\theta')$.

$\max_i g_i(\gamma\theta + (1-\gamma)\theta') \leq \max_i \{ \gamma g_i(\theta) + (1-\gamma)g_i(\theta') \}$
 $\leq \max_i \gamma g_i(\theta) + \max_i (1-\gamma)g_i(\theta')$
 $\leq \gamma \max_i g_i(\theta) + (1-\gamma) \max_i g_i(\theta')$

(3) For $d=1$. (Exercise : proof for $d \geq 1$)

$$\varphi'(\theta) = h'(g(\theta)) g'(\theta)$$

$$\varphi''(\theta) = \underbrace{h''(g(\theta))}_{\geq 0} [g'(\theta)]^2 + g''(\theta) h'(g(\theta)).$$

Use second-order condition for convexity.

ERM AS A CONVEX OPTIMIZATION PROBLEM

In some simple cases, ERM becomes a convex optimization problem.

EXAMPLE 1 $\mathcal{X} = \mathbb{R}^d$, $\mathcal{Y} = \mathbb{R}$, $\Theta \subset \mathbb{R}^d$ compact and convex set.

$$\mathcal{H}_\Theta = \{x \mapsto \theta^T x : \theta \in \Theta, x \in \mathbb{R}^d\}.$$

Then,

$$\min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n (y_i - \theta^T x_i)^2 \text{ is a convex problem.}$$

Pf: Exercise

Convexification of binary hypothesis testing:

Recall: $\hat{R}_s(f_\theta) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{f_\theta(x_i) \neq y_i\}$ for binary classification. Not convex, not continuous.

Idea: Use $h_\theta : \mathcal{X} \rightarrow \mathbb{R}$, $\theta \in \Theta$.

$$f_\theta(x) = \text{sgn}(h_\theta(x)) \text{ for } \text{sgn}(z) = \begin{cases} -1, & z < 0 \\ 0, & z = 0 \\ 1, & z > 0. \end{cases}$$

Then,

$$\{f_\theta(x_i) \neq y_i\} = \{y_i h_\theta(x_i) \leq 0\}.$$

Thus, defining

$$\ell_{0-1}(z) = \mathbb{1}\{z \leq 0\},$$

$$\hat{R}_s(f_\theta) = \frac{1}{n} \sum_{i=1}^n \ell_{0-1}(y_i h_\theta(x_i)), \quad \forall \theta \in \Theta.$$

Idea: Use convex surrogates for ℓ_{0-1} .

$$\text{Square loss : } \ell_{sq}(z) = (z-1)^2,$$

$$\text{Logistic loss : } \ell_{\text{Log}}(z) = \log(1+e^{-z}),$$

$$\text{Hinge loss : } \ell_{\text{Hinge}}(z) = \max\{0, 1-z\}.$$

Then,

PROPI For convex Θ , if $\theta \mapsto h_\theta$ is a concave function, then $\theta \mapsto \hat{R}_s(f_\theta)$ is a convex function with ℓ_{Hinge} and ℓ_{Log} .