

Optimization for Machine Learning

机器学习中的优化方法

陈程

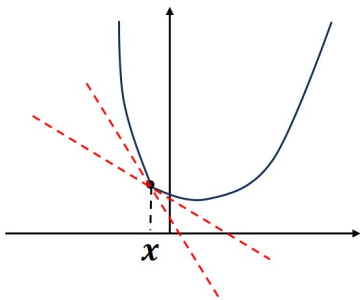
华东师范大学 软件工程学院

chchen@sei.ecnu.edu.cn

Outline

- 1 Review
- 2 Proximal Gradient Descent
- 3 Proximal Operator
- 4 Convergence Analysis

Subgradient



We say \mathbf{g} is a **subgradient** of f at the point \mathbf{x} if

$$f(\mathbf{y}) \geq \underbrace{f(\mathbf{x}) + \langle \mathbf{g}, \mathbf{y} - \mathbf{x} \rangle}_{\text{a linear under-estimate of } f}, \quad \forall \mathbf{y} \in \text{dom } f$$

The set of all subgradients of f at \mathbf{x} is called the **subdifferential** of f at \mathbf{x} , denoted by $\partial f(\mathbf{x})$.

Clean Up of Last Class

In each iteration, the (projected) subgradient descent method computes

$$\mathbf{x}_{t+1} = \mathcal{P}_C(\mathbf{x}_t - \eta_t \mathbf{g}_t),$$

where \mathbf{g}_t is **any** subgradient of f at \mathbf{x}_t .

Polyak's stepsize:

$$\eta_t = \frac{f(\mathbf{x}_t) - f^*}{\|\mathbf{g}_t\|_2^2}$$

- require to **know** f^*

Example: Finding a point in the intersection of convex sets

$\mathcal{C} = \mathcal{C}_1 \cap \dots \cap \mathcal{C}_m$ is nonempty, $\mathcal{C}_1, \dots, \mathcal{C}_m$ are closed and convex.

Find a point in \mathcal{C} by minimizing

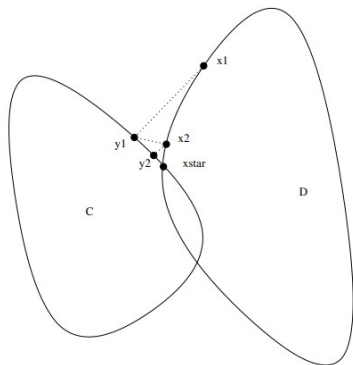
$$f(\mathbf{x}) = \max\{\text{dist}_{\mathcal{C}_1}(\mathbf{x}), \dots, \text{dist}_{\mathcal{C}_m}(\mathbf{x})\},$$

where $\text{dist}_{\mathcal{C}}(\mathbf{x}) \triangleq \min_{\mathbf{y} \in \mathcal{C}} \|\mathbf{x} - \mathbf{y}\|_2$.

With $\text{dist}_{\mathcal{C}_j}(\mathbf{x}) = f(\mathbf{x})$, a subgradient of f is

$$\mathbf{g} = \frac{\mathbf{x} - \mathcal{P}_{\mathcal{C}_j}(\mathbf{x})}{\|\mathbf{x} - \mathcal{P}_{\mathcal{C}_j}(\mathbf{x})\|_2} \in \partial \text{dist}_{\mathcal{C}_j}(\mathbf{x})$$

Example: projection onto intersection of convex sets



the subgradient method with Polyak's stepsize rule:

$$\mathbf{x}_{t+1} = \mathcal{P}_{C_j}(\mathbf{x}_t)$$

- equivalent to the famous alternating projection algorithm
- at each step, project the current point onto the farthest set

Outline

- 1 Review
- 2 Proximal Gradient Descent**
- 3 Proximal Operator
- 4 Convergence Analysis

Composite Models

$$\min_{\mathbf{x}} F(\mathbf{x}) = f(\mathbf{x}) + h(\mathbf{x})$$

- f is convex and smooth
- h is convex (may not be differentiable)
- Let $F^* = \min_{\mathbf{x}} F(\mathbf{x})$ be the optimal value

Example: ℓ_1 regularized minimization:

$$\min_{\mathbf{x}} f(\mathbf{x}) + \lambda \|\mathbf{x}\|_1$$

Review of Gradient Descent

We first revisit gradient descent

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta_t \nabla f(\mathbf{x}_t)$$



$$\mathbf{x}_{t+1} = \arg \min_{\mathbf{x}} \left\{ f(\mathbf{x}_t) + \langle \nabla f(\mathbf{x}_t), \mathbf{x} - \mathbf{x}_t \rangle + \frac{1}{2\eta_t} \|\mathbf{x} - \mathbf{x}_t\|_2^2 \right\}$$

How about projected gradient descent?

$$\mathbf{x}_{t+1} = \mathcal{P}_{\mathcal{C}}(\mathbf{x}_t - \eta_t \nabla f(\mathbf{x}_t))$$

\Leftrightarrow

$$\begin{aligned}\mathbf{x}_{t+1} &= \arg \min_{\mathbf{x}} \left\{ f(\mathbf{x}_t) + \langle \nabla f(\mathbf{x}_t), \mathbf{x} - \mathbf{x}_t \rangle + \frac{1}{2\eta_t} \|\mathbf{x} - \mathbf{x}_t\|_2^2 + \mathbb{1}_{\mathcal{C}}(\mathbf{x}) \right\} \\ &= \arg \min_{\mathbf{x}} \left\{ \frac{1}{2} \|\mathbf{x} - (\mathbf{x}_t - \eta_t \nabla f(\mathbf{x}_t))\|_2^2 + \eta_t \mathbb{1}_{\mathcal{C}}(\mathbf{x}) \right\}\end{aligned}$$

where

$$\mathbb{1}_{\mathcal{C}}(\mathbf{x}) = \begin{cases} 0, & \text{if } \mathbf{x} \in \mathcal{C} \\ +\infty, & \text{otherwise} \end{cases}$$

Proximal Operator (邻近算子)

Define the proximal operator

$$\text{prox}_h(\mathbf{x}) \triangleq \arg \min_{\mathbf{z}} \left\{ \frac{1}{2} \|\mathbf{x} - \mathbf{z}\|_2^2 + h(\mathbf{z}) \right\}$$

for any convex function h .

Then, the update of projected gradient descent is

$$\mathbf{x}_{t+1} = \text{prox}_{\eta_t \mathbb{1}_C}(\mathbf{x}_t - \eta_t \nabla f(\mathbf{x}_t))$$

Proximal Gradient Descent (邻近梯度下降法)

In each iteration, the proximal gradient descent method for composite objective function $F(\mathbf{x}) = f(\mathbf{x}) + h(\mathbf{x})$ computes

$$\mathbf{x}_{t+1} = \text{prox}_{\eta_t h}(\mathbf{x}_t - \eta_t \nabla f(\mathbf{x}_t)).$$

- alternates between gradient updates on f and proximal minimization on h
- useful if the prox can be efficiently computed

Outline

- 1 Review
- 2 Proximal Gradient Descent
- 3 Proximal Operator**
- 4 Convergence Analysis

Proximal Operator

$$\text{prox}_h(\mathbf{x}) \triangleq \arg \min_{\mathbf{z}} \left\{ \frac{1}{2} \|\mathbf{x} - \mathbf{z}\|_2^2 + h(\mathbf{z}) \right\}$$

- well-defined under very general conditions (including nonsmooth convex functions)
- can be evaluated efficiently for many widely used functions (in particular, regularizers)
- this abstraction is mathematically simple but covers many well-known optimization algorithms

Example: Indicator Functions

If $h(\mathbf{x}) = \mathbb{1}_{\mathcal{C}}$ is the “indicator” function

$$h(\mathbf{x}) = \begin{cases} 0, & \text{if } \mathbf{x} \in \mathcal{C} \\ +\infty, & \text{otherwise} \end{cases}$$

then

$$\text{prox}_h(\mathbf{x}) = \arg \min_{\mathbf{z} \in \mathcal{C}} \|\mathbf{z} - \mathbf{x}\|_2^2 \quad (\text{Euclidean projection})$$

Example: ℓ_1 Norm

If $h(\mathbf{x}) = \lambda \|\mathbf{x}\|_1$, then

$$(\text{prox}_{\lambda h}(\mathbf{x}))_i = \psi_{st}(x_i; \lambda) \quad \text{soft-thresholding}$$

where

$$\psi(x) = \begin{cases} x - \lambda, & \text{if } x > \lambda \\ x + \lambda, & \text{if } x < -\lambda \\ 0, & \text{otherwise} \end{cases}$$

Basic Rules of Proximal Operator

- **affine addition:** if $f(\mathbf{x}) = g(\mathbf{x}) + \mathbf{a}^\top \mathbf{x} + b$, then

$$\text{prox}_f(\mathbf{x}) = \text{prox}_g(\mathbf{x} - \mathbf{a})$$

- **quadratic addition:** if $f(\mathbf{x}) = g(\mathbf{x}) + \frac{\rho}{2} \|\mathbf{x} - \mathbf{a}\|_2^2$, then

$$\text{prox}_f(\mathbf{x}) = \text{prox}_{\frac{1}{1+\rho}g} \left(\frac{1}{1+\rho} \mathbf{x} - \frac{\rho}{1+\rho} \mathbf{a} \right)$$

- **scaling and translation:** if $f(\mathbf{x}) = g(a\mathbf{x} + b)$, then

$$\text{prox}_f(\mathbf{x}) = \frac{1}{a} \left(\text{prox}_{a^2g}(a\mathbf{x} + b) - b \right)$$

Basic Rules of Proximal Operator

- **norm composition:** if $f(\mathbf{x}) = g(\|\mathbf{x}\|_2)$ with $\text{dom}g = [0, \infty)$, then

$$\text{prox}_f(\mathbf{x}) = \text{prox}_g(\|\mathbf{x}\|_2) \frac{\mathbf{x}}{\|\mathbf{x}\|_2}, \quad \forall \mathbf{x} \neq \mathbf{0}$$

Nonexpansiveness of Proximal Operators

- **(firm nonexpansiveness)**

$$\langle \text{prox}_h(\mathbf{x}_1) - \text{prox}_h(\mathbf{x}_2), \mathbf{x}_1 - \mathbf{x}_2 \rangle \geq \|\text{prox}_h(\mathbf{x}_1) - \text{prox}_h(\mathbf{x}_2)\|_2^2$$

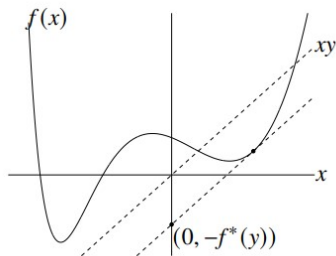
- **(nonexpansiveness)**

$$\|\text{prox}_h(\mathbf{x}_1) - \text{prox}_h(\mathbf{x}_2)\|_2 \leq \|\mathbf{x}_1 - \mathbf{x}_2\|_2$$

Conjugate Functions (共轭函数)

The **conjugate** of a function f is

$$f^*(\mathbf{y}) = \sup_{\mathbf{x} \in \text{dom } f} \{ \langle \mathbf{y}, \mathbf{x} \rangle - f(\mathbf{x}) \}$$



Conjugate Functions

Property: If f is convex and closed. Then

- $\mathbf{y} \in \partial f(\mathbf{x}) \iff \mathbf{x} \in \partial f^*(\mathbf{y})$
- $f^{**} = f$

Examples:

- **Indicator function:**

$$f(\mathbf{x}) = \mathbb{1}_{\mathcal{C}}(\mathbf{x}), \quad f^*(\mathbf{y}) = \sup_{\mathbf{x} \in \mathcal{C}} \langle \mathbf{x}, \mathbf{y} \rangle$$

- **Norm:**

$$f(\mathbf{x}) = \|\mathbf{x}\|, \quad f^*(\mathbf{y}) = \begin{cases} 0, & \|\mathbf{y}\|_* \leq 1 \\ +\infty, & \|\mathbf{y}\|_* > 1 \end{cases}$$

where $\|\mathbf{y}\|_* = \sup_{\|\mathbf{x}\| \leq 1} \langle \mathbf{x}, \mathbf{y} \rangle$ is the dual norm.

Moreau Decomposition

Suppose f is closed and convex. Then

$$\mathbf{x} = \text{prox}_f(\mathbf{x}) + \text{prox}_{f^*}(\mathbf{x})$$

Example: prox of support function

For any closed and convex set \mathcal{C} , the support function is defined as $S_{\mathcal{C}}(\mathbf{x}) = \sup_{\mathbf{z} \in \mathcal{C}} \langle \mathbf{x}, \mathbf{z} \rangle$. Then

$$\text{prox}_{S_{\mathcal{C}}}(\mathbf{x}) = \mathbf{x} - \mathcal{P}_{\mathcal{C}}(\mathbf{x})$$

Examples

- ℓ_∞ norm:

$$\text{prox}_{\|\cdot\|_\infty}(\mathbf{x}) = \mathbf{x} - \mathcal{P}_{\mathcal{B}_{\|\cdot\|_1}}(\mathbf{x})$$

where $\mathcal{B}_{\|\cdot\|_1} = \{\mathbf{z} \mid \|\mathbf{z}\|_1 \leq 1\}$ is unit ℓ_1 ball.

- **max function:** Let $g(\mathbf{x}) = \{x_1, \dots, x_n\}$, then

$$\text{prox}_g(\mathbf{x}) = \mathbf{x} - \mathcal{P}_\Delta(\mathbf{x})$$

where $\Delta = \{\mathbf{z} \in \mathbb{R}_+^n \mid \mathbf{1}^\top \mathbf{z} = 1\}$ is probability simplex.

Outline

- 1 Review
- 2 Proximal Gradient Descent
- 3 Proximal Operator
- 4 Convergence Analysis**

Convergence Analysis

Lemma

Let $\mathbf{y}^+ = \text{prox}_{\frac{1}{L}h}(\mathbf{y} - \frac{1}{L}\nabla f(\mathbf{y}))$, then

$$F(\mathbf{y}^+) - F(\mathbf{x}) \leq \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 - \frac{L}{2} \|\mathbf{x} - \mathbf{y}^+\|_2^2 - g(\mathbf{x}, \mathbf{y})$$

where $g(\mathbf{x}, \mathbf{y}) = f(\mathbf{x}) - f(\mathbf{y}) - \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle$.

- Take $\mathbf{x} = \mathbf{y} = \mathbf{x}_t$, we get $F(\mathbf{x}_{t+1}) \leq F(\mathbf{x}_t)$.
- Take $\mathbf{x} = \mathbf{x}^*$, $\mathbf{y} = \mathbf{x}_t$, we get $\|\mathbf{x}_{t+1} - \mathbf{x}^*\|_2 \leq \|\mathbf{x}_t - \mathbf{x}^*\|_2$.

Convergence for Convex Problems

Suppose f is convex and L -smooth. The proximal gradient descent with stepsize $\eta_t = 1/L$ obeys

$$F(\mathbf{x}_t) - F(\mathbf{x}^*) \leq \frac{L \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2}{2t}.$$

- Achieves better iteration complexity ($O(1/\varepsilon)$) than subgradient method ($O(1/\varepsilon^2)$).

Convergence for Strongly Convex Problems

Suppose f is μ -strongly convex and L -smooth. The proximal gradient descent with stepsize $\eta_t = 1/L$ obeys

$$\|\mathbf{x}_t - \mathbf{x}^*\|_2^2 \leq \left(1 - \frac{\mu}{L}\right)^t \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2.$$

- Achieves linear convergence $O(\kappa \log \frac{1}{\varepsilon})$.

Summary

condition	stepsize	convergence rate	iteration complexity
convex & smooth	$\eta_t = \frac{1}{L}$	$O\left(\frac{1}{t}\right)$	$O\left(\frac{1}{\varepsilon}\right)$
strongly convex & smooth	$\eta_t = \frac{1}{L}$	$O\left(\left(1 - \frac{1}{\kappa}\right)^t\right)$	$O\left(\kappa \log \frac{1}{\varepsilon}\right)$

Table: Convergence Properties of Proximal Gradient Descent

	stepsize	convergence rate	iteration complexity
convex & smooth	$\eta_t \approx \frac{1}{\sqrt{t}}$	$O\left(\frac{1}{\sqrt{t}}\right)$	$O\left(\frac{1}{\varepsilon^2}\right)$
strongly convex & smooth	$\eta_t \approx \frac{1}{t}$	$O\left(\frac{1}{t}\right)$	$O\left(\frac{1}{\varepsilon}\right)$

Table: Convergence Properties of Subgradient Descent

Questions

