

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

- 1 Unconstrained Optimization
- 2 Quadratic Minimization Problems
- 3 Regularity Conditions

- 1 Unconstrained Optimization
- 2 Quadratic Minimization Problems
- 3 Regularity Conditions

Differentiable Unconstrained Optimization

Suppose the objective function (or loss function) f is differentiable. The unconstrained optimization problem is:

$$\begin{aligned} \min_{\mathbf{x}} f(\mathbf{x}) \\ \text{s.t. } \mathbf{x} \in \mathbb{R}^d \end{aligned}$$

Optimal Condition (最优性条件)

Suppose f is differentiable and convex. A point \mathbf{x}^* is optimal if and only if

$$\nabla f(\mathbf{x}^*) = 0.$$

Strict convex function has **unique** optimal solution.

Iterative Descent Methods

Start with a point \mathbf{x}_0 and construct a sequence $\{\mathbf{x}_t\}$ s.t.,

$$f(\mathbf{x}_{t+1}) < f(\mathbf{x}_t). \quad t = 0, 1, \dots$$

We call \mathbf{d} is a **descent direction** at \mathbf{x} if

$$f'(\mathbf{x}; \mathbf{d}) \triangleq \underbrace{\lim_{t \rightarrow 0} \frac{f(\mathbf{x} + t\mathbf{d}) - f(\mathbf{x})}{t}}_{\text{directional derivative}} = \nabla f(\mathbf{x})^\top \mathbf{d} < 0.$$

Iterative Descent Methods

- Start with a point \mathbf{x}_0 ;
- In each iteration, search in descent direction

$$\mathbf{x}_{t+1} = \mathbf{x}_t + \eta_t \mathbf{d}_t$$

where \mathbf{d}_t is the descent direction at \mathbf{x}_t and η_t is the stepsize.

How to Find a Descent Direction?

By Cauchy-Schwarz inequality,

$$\min_{\|\mathbf{d}\|_2 \leq 1} f'(\mathbf{x}; \mathbf{d}) = \min_{\|\mathbf{d}\|_2 \leq 1} \nabla f(\mathbf{x})^\top \mathbf{d} = -\|\nabla f(\mathbf{x})\|_2$$

$f'(\mathbf{x}; \mathbf{d})$ achieve minimum when $\mathbf{d} = -\nabla f(\mathbf{x})$.

Gradient Descent (梯度下降法)

One of the most important descent methods: gradient descent

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

- descent direction: $\mathbf{d}_t = -\nabla f(\mathbf{x}_t)$
- traced to Augustin Louis Cauchy '1847
- First-order Taylor approximation: $f(\mathbf{x}) \approx f(\mathbf{x}_t) + \langle \nabla f(\mathbf{x}_t), \mathbf{x} - \mathbf{x}_t \rangle$

Outline

- 1 Unconstrained Optimization
- 2 Quadratic Minimization Problems
- 3 Regularity Conditions

Quadratic Minimization

We begin with the quadratic objective function:

$$\min_{\mathbf{x}} f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} - \mathbf{b}^\top \mathbf{x},$$

for some $d \times d$ symmetric matrix $\mathbf{Q} \succ 0$.

- The gradient is $\nabla f(\mathbf{x}) = \mathbf{Q} \mathbf{x} - \mathbf{b}$.
- The unique optimal solution is $\mathbf{x}^* = \mathbf{Q}^{-1} \mathbf{b}$.
- $\lambda_1(\mathbf{Q}) \mathbf{I} \succeq \mathbf{Q} \succeq \lambda_d(\mathbf{Q}) \mathbf{I}$, where $\lambda_1(\mathbf{Q})$ and $\lambda_d(\mathbf{Q})$ are largest and smallest eigenvalues of \mathbf{Q} respectively.

How to Find a Good Stepsize?

According to the GD update rule,

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

$$\Rightarrow \|\mathbf{x}_{t+1} - \mathbf{x}^*\|_2 \leq \|\mathbf{I} - \eta_t \mathbf{Q}\|_2 \|\mathbf{x}_t - \mathbf{x}^*\|_2$$

We observe that

$$\begin{aligned} \|\mathbf{I} - \eta_t \mathbf{Q}\|_2 &= \underbrace{\max\{|1 - \eta_t \lambda_1(\mathbf{Q})|, |1 - \eta_t \lambda_d(\mathbf{Q})|\}}_{\text{optimal choice is } \eta_t = \frac{2}{\lambda_1(\mathbf{Q}) + \lambda_d(\mathbf{Q})}} \\ &= \frac{\lambda_1(\mathbf{Q}) - \lambda_d(\mathbf{Q})}{\lambda_1(\mathbf{Q}) + \lambda_d(\mathbf{Q})} \end{aligned}$$

Convergence for Constant Stepsize

If $\eta_t = \eta = \frac{2}{\lambda_1(\mathbf{Q}) + \lambda_d(\mathbf{Q})}$, then

$$\|\mathbf{x}_t - \mathbf{x}^*\|_2 \leq \left(\frac{\lambda_1(\mathbf{Q}) - \lambda_d(\mathbf{Q})}{\lambda_1(\mathbf{Q}) + \lambda_d(\mathbf{Q})} \right)^t \|\mathbf{x}_0 - \mathbf{x}^*\|_2.$$

The stepsize $\eta_t = \eta = \frac{2}{\lambda_1(\mathbf{Q}) + \lambda_d(\mathbf{Q})}$ relies on the eigenvalues of \mathbf{Q} , which requires preliminary experimentation.

Outline

- 1 Unconstrained Optimization
- 2 Quadratic Minimization Problems
- 3 Regularity Conditions

Generalization

Let's now generalize quadratic minimization to a broader class of problems

$$\min_{\mathbf{x}} f(\mathbf{x})$$

where

$$\mu \mathbf{I} \preceq \nabla^2 f(\mathbf{x}) \preceq L \mathbf{I}.$$

Smoothness (光滑性)

We say that a function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is G -Lipschitz continuous if for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, we have

$$|f(\mathbf{x}) - f(\mathbf{y})| \leq G \|\mathbf{x} - \mathbf{y}\|_2.$$

We say a differentiable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is L -smooth if it has L -Lipschitz continuous gradient. That is, for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$, we have

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 \leq L \|\mathbf{x} - \mathbf{y}\|_2.$$

Which of following functions are smooth?

- $f(x) = x^4$;
- $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top \mathbf{Q}\mathbf{x} - \mathbf{b}^\top \mathbf{x}$ with $\mathbf{Q} \succeq 0$;
- $f(x) = \sin x$.

Equivalent First-Order Characterizations of Smoothness

Let $f : \mathbb{R}^d \leftarrow \mathbb{R}$ be a **convex** and differentiable function. Then the following properties are equivalent characterizations of L -smoothness of f :

- 1 $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 \leq L\|\mathbf{x} - \mathbf{y}\|_2, \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d;$
- 2 $\langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \leq L\|\mathbf{x} - \mathbf{y}\|_2^2, \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d;$
- 3 $f(\mathbf{y}) \leq \underbrace{f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle}_{\text{first-order Taylor expansion}} + \frac{L}{2}\|\mathbf{x} - \mathbf{y}\|_2^2, \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d;$
- 4 $f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{1}{2L}\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2^2, \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d;$
- 5 $\langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \geq \frac{1}{L}\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2^2, \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d;$

Think: which characterizations do not hold if f is not convex?

Equivalent Second-Order Characterization of Smoothness

We say a differentiable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is **L -smooth** if

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 \leq L \|\mathbf{x} - \mathbf{y}\|_2.$$

Second-Order Characterization:

Let $f : \mathbb{R}^d \leftarrow \mathbb{R}$ be a twice differentiable function. Then the following property is an equivalent characterization of L -smoothness of f :

$$-L\mathbf{I} \preceq \nabla^2 f(\mathbf{x}) \preceq L\mathbf{I}.$$

Strongly Convexity (强凸性)

We say f is μ -strongly convex if the function

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

is convex for some $\mu > 0$.

Equivalent First-Order Characterizations of Strong Convexity

Let $f : \mathbb{R}^d \leftarrow \mathbb{R}$ be a convex and differentiable function. Then the following properties are equivalent characterizations of μ -strong convexity of f :

- 1 $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 \geq \mu\|\mathbf{x} - \mathbf{y}\|_2, \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d;$
- 2 $\langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \geq \mu\|\mathbf{x} - \mathbf{y}\|_2^2, \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d;$
- 3 $f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{\mu}{2}\|\mathbf{x} - \mathbf{y}\|_2^2, \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d;$
- 4 $f(\mathbf{y}) \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{1}{2\mu}\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2^2, \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d;$
- 5 $\langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \leq \frac{1}{\mu}\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2^2, \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d;$

Strongly convex functions are strictly convex.

Equivalent Second-Order Characterization of Strongly Convexity

Second-Order Characterization:

Let $f : \mathbb{R}^d \leftarrow \mathbb{R}$ be a twice differentiable function. Then the following property is an equivalent characterization of μ -strongly convex of f :

$$\nabla^2 f(\mathbf{x}) \succeq \mu \mathbf{I}.$$

Strongly Convex and Smooth Functions

Let f be L -smooth and μ -strongly convex. Then we have

$$\mu \mathbf{I} \preceq \nabla^2 f(\mathbf{x}) \preceq L \mathbf{I}.$$

Let $\kappa \triangleq \frac{L}{\mu}$ be the **condition number**.

Convergence Rate of Strongly Convex and Smooth Problems

Let f be L -smooth and μ -strongly convex. If $\eta_t = \eta = \frac{2}{\mu+L}$, then

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

Iteration complexity: To achieve ϵ -accuracy, we require $\frac{\log(\|\mathbf{x}_0 - \mathbf{x}^*\|_2 / \epsilon)}{\log(\frac{\kappa+1}{\kappa-1})}$ number of iterations.

Dimension-free: The iteration complexity is independent of problem size d if κ does not depend on d .