

Optimization Theory

Lecture 05

Fudan University

luoluo@fudan.edu.cn

- 1 Second-Order Characterization
- 2 Examples and Applications
- 3 Black Box Model

1 Second-Order Characterization

2 Examples and Applications

3 Black Box Model

Second-Order Characterization

Theorem (Smoothness and Convexity)

Let $f(\cdot)$ be a twice differentiable function defined on \mathbb{R}^d

- ① It is L -smooth if and only if $-L\mathbf{I} \preceq \nabla^2 f(\mathbf{x}) \preceq L\mathbf{I}$ for all $\mathbf{x} \in \mathbb{R}^d$.
- ② It is convex if and only if $\nabla^2 f(\mathbf{x}) \succeq \mathbf{0}$ for all $\mathbf{x} \in \mathbb{R}^d$.
- ③ It is μ -strongly-convex if and only if $\nabla^2 f(\mathbf{x}) \succeq \mu\mathbf{I}$ for all $\mathbf{x} \in \mathbb{R}^d$.

Sometimes, we say $f(\cdot)$ is ℓ -weakly convex if the function

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

is convex for some $\ell > 0$.

Theorem

Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a twice differentiable function. Suppose that $\nabla^2 f(\cdot)$ is continuous in an open neighborhood of $\mathbf{x}^* \in \mathbb{R}^d$.

① If \mathbf{x}^* is a local minimizer of $f(\cdot)$, then it holds that

$$\nabla f(\mathbf{x}^*) = \mathbf{0} \quad \text{and} \quad \nabla^2 f(\mathbf{x}^*) \succeq \mathbf{0}.$$

② If it holds that

$$\nabla f(\mathbf{x}^*) = \mathbf{0} \quad \text{and} \quad \nabla^2 f(\mathbf{x}^*) \succ \mathbf{0},$$

then the point \mathbf{x}^* is a strict local minimizer of $f(\cdot)$.

1 Second-Order Characterization

2 Examples and Applications

3 Black Box Model

Examples

- ① For unconstrained quadratic problem

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) \triangleq \frac{1}{2} \mathbf{x}^\top \mathbf{A} \mathbf{x} - \mathbf{b}^\top \mathbf{x},$$

where $\mathbf{A} \in \mathbb{R}^{d \times d}$. We have $\nabla^2 f(\mathbf{x}) = \mathbf{A}$.

- ② For regularized generalized linear model

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) \triangleq \frac{1}{n} \sum_{i=1}^n \phi_i(\mathbf{a}_i^\top \mathbf{x}) + \frac{\lambda}{2} \|\mathbf{x}\|_2^2.$$

where $\phi_i : \mathbb{R}^d \rightarrow \mathbb{R}$ is twice differentiable. We have

$$\nabla f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \phi'_i(\mathbf{a}_i^\top \mathbf{x}) \mathbf{a}_i + \lambda \mathbf{x}$$

and

$$\nabla^2 f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \phi''_i(\mathbf{a}_i^\top \mathbf{x}) \mathbf{a}_i \mathbf{a}_i^\top + \lambda \mathbf{I}.$$

Applications in Matrix Approximation

Given a symmetric positive-definite matrix $\mathbf{K} \in \mathbb{R}^{d \times d}$ and we sample a subset of columns $\mathbf{C} \in \mathbb{R}^{d \times m}$, where $m < d$.

We want to establish the estimator of \mathbf{K} by the formulation

$$\min_{\mathbf{U} \in \mathbb{R}^{m \times m}, \delta \in \mathbb{R}} f(\mathbf{U}, \delta) \triangleq \left\| \mathbf{K} - (\mathbf{CUC}^\top + \delta \mathbf{I}_d) \right\|_F^2.$$

It has global solution

$$\mathbf{U}^{\text{ss}} = \mathbf{C}^\dagger \mathbf{K} (\mathbf{C}^\dagger)^\top - \delta^{\text{ss}} (\mathbf{C}^\top \mathbf{C})^\dagger$$

and

$$\delta^{\text{ss}} = \frac{1}{d - m} \left(\text{tr}(\mathbf{K}) - \text{tr}(\mathbf{C}^\dagger \mathbf{K} \mathbf{C}) \right).$$

Applications in Matrix Approximation

We can show that

$$\mathbf{C}\mathbf{U}^{\text{ss}}\mathbf{C}^\top + \delta^{\text{ss}}\mathbf{I}_d \succ \mathbf{0}$$

and

$$\begin{aligned} & (\mathbf{Q}\mathbf{U}^{\text{ss}}\mathbf{Q}^\top + \delta^{\text{ss}}\mathbf{I}_d)^{-1} \\ &= (\delta^{\text{ss}})^{-1}\mathbf{I}_n - (\delta^{\text{ss}})^{-2}\mathbf{Q}(\mathbf{I}_m + (\delta^{\text{ss}})^{-1}\mathbf{U}^{\text{ss}})^{-1}\mathbf{U}^{\text{ss}}\mathbf{Q}^\top. \end{aligned}$$

is well-defined, where $\mathbf{Q} \in \mathbb{R}^{d \times m}$ is the orthogonal basis of $\mathbf{C} \in \mathbb{R}^{d \times m}$.

Outline

1 Second-Order Characterization

2 Examples and Applications

3 Black Box Model

Convergence Criteria

For the unconstrained convex optimization problem

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}),$$

the convergence of an algorithm can be measured by the following in metrics:

- 1 Convergence in parameter (suppose there exists optimal solution \mathbf{x}^*), where we measure the distance

$$\|\mathbf{x}_t - \mathbf{x}^*\|_2.$$

- 2 Convergence of objective value, measured by objective suboptimality

$$f(\mathbf{x}_t) - \inf_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}).$$

- 3 Convergence of gradient

$$\|\nabla f(\mathbf{x}_t)\|_2.$$

Convergence Criteria

If $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is smooth and convex and has an optimal solution \mathbf{x}^* , then

$$f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \langle \nabla f(\mathbf{x}^*), \mathbf{x}_t - \mathbf{x}^* \rangle + \frac{L}{2} \|\mathbf{x}_t - \mathbf{x}^*\|_2^2 = \frac{L}{2} \|\mathbf{x}_t - \mathbf{x}^*\|_2^2,$$

and

$$\|\nabla f(\mathbf{x}_t)\|_2 = \|\nabla f(\mathbf{x}_t) - \nabla f(\mathbf{x}^*)\|_2 \leq L \|\mathbf{x}_t - \mathbf{x}^*\|_2,$$

which implies convergence in parameter implies convergence in objective value and gradient.

The reverse directions may not hold if the objective is not strongly-convex.

Black Box Model

Local black box:

- ① The only information available for the numerical scheme is the answer of the oracle.
- ② The oracle is local.

Different types of oracle:

- ① Zero-order oracle: returns the function value $f(\mathbf{x})$.
- ② First-order oracle: returns the function value $f(\mathbf{x})$ and the gradient $\nabla f(\mathbf{x})$.
- ③ Second-order oracle: returns $f(\mathbf{x})$, $\nabla f(\mathbf{x})$, and the Hessian $\nabla^2 f(\mathbf{x})$.

Black Box Model

There are two participants in the black box model: a learner and an oracle.

- ① The learner has
 - infinite computational power,
 - knowledge of the function class to which f belongs,
 - knowledge of the domain.
- ② The oracle has specific knowledge of the function.

The key question:

How many queries to the oracles are necessary and sufficient to find an ϵ -approximate solution?

We will study this question from two perspectives:

- ① Upper bound: Designing algorithms.
- ② Lower bound: Information theoretic reasoning.

The strength of the black-box model:

- ① It will allow us to derive a complete theory of optimization.
- ② We will obtain matching upper and lower bounds on the oracle complexity for various sub-classes of interesting functions.

The weakness of the black-box model:

- ① It does not limit our computational resources.
- ② The side information of the algorithm is ignored.