

# Numerical Optimization

Graduate Course

## Unconstrained Smooth Optimization

Part I: Basic theories and differentials

Wen Huang

School of Mathematical Sciences  
Xiamen University

Compiled on February 14, 2022

# Preliminaries and Basic Theories

# Line Search Methods

A representative line search method

---

## A descent method

**Input:** Initial iterate  $x_0$ ;

Initial descent direction  $p_0$  and set  $k \leftarrow 0$ ;

**while** not accurate enough **do**

Set  $x_{k+1} \leftarrow x_k + \alpha_k p_k$  with an appropriate step size  $\alpha_k$ ;

Find a descent direction  $p_{k+1}$ ;

$k \leftarrow k + 1$ ;

**end while**

---

**What is the definition of a descent direction?**

# Line Search Methods

## Descent directions

### Definition 1

Let  $f \in C^1$  and  $p \in \mathbb{R}^n$ .  $p$  is a descent direction of  $f$  at  $x$  if for all sufficiently small  $\alpha > 0$

$$f(x + \alpha p) < f(x)$$

# Line Search Methods

## Descent directions

### Theorem 2

Let  $f \in C^1$  and  $p \in \mathbb{R}^n$ . If

$$-p^T \nabla f(x_k) = \|p\|_2 \|\nabla f(x_k)\|_2 \cos \theta_k > 0$$

then  $p$  is a descent direction at  $x$ .

Note that a direction  $p$  is descent for  $f \in C^1$  if the angle between  $p$  and  $-\nabla f(x_k)$  is acute.

To prove this theorem, we need the Taylor's theorem of multiple variables.

# Taylor's Theorem

## Theorem 3 (Taylor's Theorem)

Let  $f$  be a real-value function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ .

- If  $f \in C^1$  and  $p \in \mathbb{R}^n$  then for some  $0 < \tau < 1$

$$f(x + p) = f(x) + \nabla f(x + \tau p)^T p$$

- If  $f \in C^2$  and  $p \in \mathbb{R}^n$  then for some  $0 < \tau < 1$

$$f(x + p) = f(x) + \nabla f(x)^T p + 0.5 p^T \nabla^2 f(x + \tau p) p$$

and

$$\nabla f(x + p) = \nabla f(x) + \int_0^1 \nabla^2 f(x + \tau p) p \, d\tau$$

The proofs are omitted here.

# Minimizers

## Definition 4

The point  $x^* \in \mathbb{R}^n$  is a global minimizer if  $f(x^*) \leq f(x)$  for all  $x \in \mathbb{R}^n$ .

## Definition 5

The point  $x^* \in \mathbb{R}^n$  is a local minimizer if  $f(x^*) \leq f(x)$  for all  $x \in \mathcal{N}_{x^*}$  where  $\mathcal{N}_{x^*}$  is a neighborhood of  $x^*$ . Further,  $x^*$  is

- a strict local minimizer if  $f(x^*) < f(x)$  for all  $x \in \mathcal{N}_{x^*}$ ; and
- an isolated local minimizer if  $x^*$  is the only local minimizer in  $\mathcal{N}_{x^*}$ .

## Definition 6

The point  $x^* \in \mathbb{R}^n$  is a stationary point if  $\nabla f(x^*) = 0$ .

# Optimality Conditions

## First order necessary condition

### Theorem 7 (First order necessary condition)

*Suppose  $f \in C^1$  in a neighborhood of  $x^*$ . If  $x^*$  is a local minimizer then  $\nabla f(x^*) = 0$ , i.e.,  $x^*$  is a stationary point.*



# Optimality Conditions

## Second order necessary condition

### Theorem 8 (Second order necessary condition)

*Suppose  $f \in C^2$  in a neighborhood of  $x^*$ . If  $x^*$  is a local minimizer then  $\nabla f(x^*) = 0$  and  $\nabla^2 f(x^*)$  is positive semidefinite.*

# Optimality Conditions

## Second order sufficient condition

### Theorem 9 (Second order sufficient condition)

*Suppose  $f \in C^2$  in a neighborhood of  $x^*$ . If  $\nabla f(x^*) = 0$  and  $\nabla^2 f(x^*)$  is positive definite then  $x^*$  is a strict local minimizer.*

# Optimality Conditions

## Sufficient condition of global minimizer

### Definition 10

A function  $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$  is convex if for any two points,  $x$  and  $y$ , in the domain we have for any  $0 \leq \alpha \leq 1$

$$f(\alpha x + (1 - \alpha)y) \leq \alpha f(x) + (1 - \alpha)f(y)$$

Also,  $f(x)$  is concave if  $-f(x)$  is convex.

### Theorem 11 (Sufficient condition of global minimizer)

*If  $f$  is convex then any local minimizer is a global minimizer. If, in addition,  $f \in C^1$  then any stationary point is a global minimizer.*

# Rates of Convergence

Desire convergence properties

- Global convergence (to a stationary point from any initial point)
  - Fast convergence speed
    - $x_k \rightarrow x_*$ , let  $e_k = \|x_k - x_*\|$ ;
    - $f_k \rightarrow f_*$ , let  $e_k = |f_k - f_*|$ ;
- 

- Q-sublinear  $\frac{e_{k+1}}{e_k} \rightarrow 1$
- Q-linear  $\limsup \frac{e_{k+1}}{e_k} < \delta < 1$
- Q-superlinear  $\frac{e_{k+1}}{e_k} \rightarrow 0$
- Q-quadratic  $\limsup \frac{e_{k+1}}{e_k^2} < C$

# Rates of Convergence

Desire convergence properties

- Global convergence (to a stationary point from any initial point)
  - Fast convergence speed
    - $x_k \rightarrow x_*$ , let  $e_k = \|x_k - x_*\|$ ;
    - $f_k \rightarrow f_*$ , let  $e_k = |f_k - f_*|$ ;
- 

- Q-sublinear  $\frac{e_{k+1}}{e_k} \rightarrow 1$
- Q-linear  $\limsup \frac{e_{k+1}}{e_k} < \delta < 1$
- Q-superlinear  $\frac{e_{k+1}}{e_k} \rightarrow 0$
- Q-quadratic  $\limsup \frac{e_{k+1}}{e_k^2} < C$

Suppose  $e_k \leq \epsilon_k$

- $\epsilon_k$  Q-sublinear  $\Rightarrow$  R-sublinear
- $\epsilon_k$  Q-linear  $\Rightarrow$  R-linear
- $\epsilon_k$  Q-superlinear  $\Rightarrow$  R-superlinear
- $\epsilon_k$  Q-quadratic  $\Rightarrow$  R-Quadratic

# Gradient and Hessian

# Gradient and Hessian

## Definition

### Definition 12

If  $f : \mathbb{R}^n \rightarrow \mathbb{R} : x \mapsto f(x)$  then

- the gradient of  $f \in C^1$ , denoted,  $\nabla f(x) \in \mathbb{R}^n$ , is the vector with  $i$ -th element,

$$\frac{\partial f}{\partial x_i}(x)$$

- the Hessian of  $f \in C^2$  denoted,  $\nabla^2 f(x) \in \mathbb{R}^{n \times n}$ , is the symmetric matrix with  $i, j$ -element,

$$\frac{\partial^2 f}{\partial x_i \partial x_j}(x)$$

# Gradient and Hessian

## Examples

- Example:  $f(x) = \frac{1}{2}x^T A x$ , where  $A = A^T$



# Gradient and Hessian

## Examples

- Example:  $f(x) = \frac{1}{2}x^T Ax$ , where  $A = A^T$ 
  - By definition

# Gradient and Hessian

## Examples

- Example:  $f(x) = \frac{1}{2}x^T Ax$ , where  $A = A^T$ 
  - By definition
  - By directional derivative
    - $Df(x)[d] = \nabla f(x)^T d, \forall d \in \mathbb{R}^n$ ;
    - $D(\nabla f(x))[d] = \nabla^2 f(x)d, \forall d \in \mathbb{R}^n$ ;

# Derivation formulas

If  $F$  is linear, then

$$DF(x)[z] = F(z).$$

Chain rule: If  $\text{range}(F) \subseteq \text{dom}(G)$ , then

$$D(G \circ F)(x)[z] = DG(F(x))[DF(x)[z]].$$

Product rule: Let  $\bullet$  denote a bilinear operator on the ranges of  $F$  and  $G$ , then

$$D(F \bullet G)(x)[z] = (DF(x)[z]) \bullet G(x) + F(x) \bullet (DG(x)[z]).$$

## Example 13

Compute the gradient and the action of Hessian of the function  
 $f(x) = \frac{1}{2} \| (xx^T - I_n)A \|_F^2$ .

# Gradient and Hessian

## Matrix function

$$\min_{x \in \mathcal{E}} f(x)$$

- In previous examples,  $\mathcal{E} = \mathbb{R}^n$ ;
- What if  $\mathcal{E}$  is the set of real matrices  $\mathbb{R}^{m \times n}$ , the set of complex matrices  $\mathbb{C}^{m \times n}$ , or even real/complex tensors;

# Gradient and Hessian

## Matrix function

Function on matrices:  $f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$

- Gradient:

$$\nabla f(X) = \begin{bmatrix} \frac{\partial f}{\partial x_{11}} & \frac{\partial f}{\partial x_{12}} & \cdots & \frac{\partial f}{\partial x_{1n}} \\ \frac{\partial f}{\partial x_{21}} & \frac{\partial f}{\partial x_{22}} & \cdots & \frac{\partial f}{\partial x_{2n}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial x_{m1}} & \frac{\partial f}{\partial x_{m2}} & \cdots & \frac{\partial f}{\partial x_{mn}} \end{bmatrix}$$

- $Df(X)[V] = \text{trace}(\nabla f(X)^T V), \forall V \in \mathbb{R}^{m \times n};$

# Gradient and Hessian

## Matrix function

### Example 14

Compute the gradient of

$$f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R} : X \mapsto \log \det(I_m + XX^T).$$

(Hint:  $D(\log \det(Y))[V] = \text{trace}(Y^{-1}V)$  for a symmetric positive definite matrix  $Y$  and symmetric matrix  $V$ .)

# Gradient and Hessian

## Matrix function

- Action of the Hessian

$$\nabla^2 f(X)[V] = D(\nabla f(X))[V];$$

### Example 15

Compute the action of the Hessian of

$$f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R} : X \mapsto \log \det(I_m + XX^T).$$



# Gradient and Hessian

## Matrix function

- No matter whether the function is a vector function or a real matrix function, we have

$$Df(x)[v] = \sum_{i_1, i_2, \dots, i_s} (\nabla f(x))_{i_1, i_2, \dots, i_s} v_{i_1, i_2, \dots, i_s},$$

where  $s = 1$  for a vector and  $s = 2$  for a matrix;

- Such idea can be extent to tensor or complex numbers or product of multiple spaces;

# Gradient and Hessian

## Product

### Example 16

Compute the gradient of

$$f : \mathbb{R}^{m \times k} \times \mathbb{R}^{n \times k} \rightarrow \mathbb{R} : (X, Y) \mapsto \text{trace}(X^T A Y),$$

where  $A \in \mathbb{R}^{m \times n}$ ;

# References I