

Numerical Optimization

Unit 2: Multivariable optimization problems

Che-Rung Lee

Scribe: 張雅芳

February 28, 2011

Partial derivative of a two variable function

- Given a two variable function $f(x_1, x_2)$.
- The partial derivative of f with respect to x_i is

$$\begin{cases} \frac{\partial f}{\partial x_1} = \lim_{h \rightarrow 0} \frac{f(x_1 + h, x_2) - f(x_1, x_2)}{h} \\ \frac{\partial f}{\partial x_2} = \lim_{h \rightarrow 0} \frac{f(x_1, x_2 + h) - f(x_1, x_2)}{h} \end{cases}$$

- The meaning of partial derivative: let $F(x_1) = f(x_1, v)$ and $G(x_2) = f(u, x_2)$,

$$\frac{\partial f}{\partial x_1}(x_1, v) = F'(x_1).$$

$$\frac{\partial f}{\partial x_2}(u, x_2) = G'(x_2).$$

Definition

The gradient of a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a **vector** in \mathbb{R}^n defined as

$$\vec{g} = \nabla f(\vec{x}) = \begin{pmatrix} \partial f / \partial x_1 \\ \vdots \\ \partial f / \partial x_n \end{pmatrix}, \text{ where } \vec{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$$

Definition

The directional derivative of a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ in the direction \vec{p} is defined as

$$D(f(\vec{x}), \vec{p}) = \lim_{h \rightarrow 0} \frac{f(\vec{x} + h\vec{p}) - f(x)}{h}$$

Remark

If $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is continuously differentiable in a neighborhood of \vec{x} ,

$$D(f(\vec{x}), \vec{p}) = \nabla f(x)^T \vec{p},$$

for any vector \vec{p} .

The descent directions

- A direction \vec{p} is called a **descent direction** of $f(\vec{x})$ at \vec{x} if $D(f(\vec{x}_0), \vec{p}) < 0$.
- If f is smooth enough, \vec{p} is a descent direction if $f(\vec{x}_0)^T \vec{p} < 0$.
- Which direction \vec{p} makes $f(\vec{x}_0 + \vec{p})$ decreasing most?
 - Mean Value theorem

$$f(\vec{x}_0 + \vec{p}) = f(\vec{x}_0) + \nabla f(\vec{x}_0 + \alpha \vec{p})^T \vec{p}$$

- $\vec{p} = -\nabla f(\vec{x}_0)$ is called the steepest descent direction of $f(x)$ at x_0 .

$$\begin{aligned} f(\vec{x}_0 + \vec{p}) &= f(\vec{x}_0) + \nabla f(\vec{x}_0 + \alpha \vec{p})^T \vec{p} \\ &\approx f(\vec{x}_0) - \nabla f(\vec{x}_0)^T \nabla f(\vec{x}_0) \end{aligned}$$

The steepest descent algorithm

The steepest descent algorithm

For $k = 1, 2, \dots$ until convergence

Compute $\vec{p}_k = -\nabla f(x_k)$

Find $\alpha_k \in (0, 1)$ s.t, $F(\alpha_k) = f(\vec{x}_k + \alpha_k \vec{p}_k)$ is minimized.

$x_{k+1} = \vec{x}_k + \alpha_k \vec{p}_k$

- You can use any single variable optimization techniques to compute α_k .
- If $F(\alpha_k) = f(\vec{x}_k + \alpha_k \vec{p}_k)$ is a quadratic function, α_k has a theoretical formula. (will be derived in next slides.)
- If $F(\alpha_k) = f(\vec{x}_k + \alpha_k \vec{p}_k)$ is more than a quadratic function, we may approximate it by a quadratic model and use the formula to solve α_k .
- Higher order polynomial approximation will be mentioned in the line search algorithm.

Quadratic model

- If $f(\vec{x})$ is a quadratic function, we can write it as

$$f(x, y) = ax^2 + bxy + cy^2 + dx + ey + f(0, 0).$$

- If f is smooth, the derivatives of f are

$$\frac{\partial f}{\partial x} = 2ax + by + d, \quad \frac{\partial f}{\partial y} = 2cy + bx + e$$

$$\frac{\partial^2 f}{\partial x^2} = 2a, \quad \frac{\partial^2 f}{\partial y^2} = 2c, \quad \frac{\partial^2 f}{\partial x \partial y} = \frac{\partial^2 f}{\partial y \partial x} = b.$$

- Let $\vec{x} = \begin{pmatrix} x \\ y \end{pmatrix}$, $f(\vec{x})$ can be expressed as

$$f(\vec{x}) = \frac{1}{2} \vec{x}^T \begin{pmatrix} 2a & b \\ b & 2c \end{pmatrix} \vec{x} + \vec{x}^T \begin{pmatrix} d \\ e \end{pmatrix} + f(\vec{0}).$$

Gradient and Hessian

- The gradient of f , as defined before, is

$$g(\vec{x}) = \nabla f(\vec{x}) = \begin{pmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{pmatrix} = \begin{pmatrix} 2a & b \\ b & 2c \end{pmatrix} \vec{x} + \begin{pmatrix} d \\ e \end{pmatrix}$$

- The second derivative, which is a matrix called **Hessian**, is

$$\nabla^2 f(\vec{x}) = H(\vec{x}) = \begin{pmatrix} \frac{\partial^2 f}{\partial x^2} & \frac{\partial^2 f}{\partial x \partial y} \\ \frac{\partial^2 f}{\partial y \partial x} & \frac{\partial^2 f}{\partial y^2} \end{pmatrix} = \begin{pmatrix} 2a & b \\ b & 2c \end{pmatrix}$$

- Therefore, $f(\vec{x}) = 1/2 \vec{x}^T H(\vec{0}) \vec{x} + g(\vec{0})^T \vec{x} + f(\vec{0})$,

$$\nabla f(\vec{x}) = H\vec{x} + \vec{g}, \text{ and } \nabla^2 f = H$$

- In the following lectures, we assume H is symmetric. Thus, $H = H^T$.

Optimal α_k for quadratic model

- We denote $H_k = H(\vec{x}_k)$, $\vec{g}_k = \vec{g}(\vec{x}_k)$, and $f_k = f(\vec{x}_k)$.
- Also, $H = H(\vec{0})$, $\vec{g} = \vec{g}(\vec{0})$, and $f = f(\vec{0})$.

$$\begin{aligned}F(\alpha) &= f(\vec{x}_k + \alpha\vec{p}_k) \\&= \frac{1}{2}(\vec{x}_k + \alpha\vec{p}_k)^T H(\vec{x}_k + \alpha\vec{p}_k) + \vec{g}^T(\vec{x}_k + \alpha\vec{p}_k) + f(\vec{0}) \\&= \frac{1}{2}\vec{x}_k^T H\vec{x}_k + \vec{g}^T \vec{x}_k + f(\vec{0}) + \alpha(H\vec{x}_k + \vec{g})^T \vec{p}_k + \frac{\alpha^2}{2}\vec{p}_k^T H\vec{p}_k \\&= f_k + \alpha\vec{g}_k^T \vec{p}_k + \frac{\alpha^2}{2}\vec{p}_k^T H\vec{p}_k \\F'(\alpha) &= \vec{g}_k^T \vec{p}_k + \alpha\vec{p}_k^T H\vec{p}_k\end{aligned}$$

The optimal solution of α_k is at $F'(\alpha) = 0$, which is $\alpha_k = \frac{-\vec{g}_k^T \vec{p}_k}{\vec{p}_k^T H\vec{p}_k}$

Theorem (Necessary and sufficient condition of optimality)

- Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be continuously differentiable in D . If $\vec{x}^* \in D$ is a local minimizer, $\nabla f(\vec{x}^*) = 0$ and $\nabla^2 f(\vec{x}^*)$ is **positive semidefinite**.
- If $\nabla f(\vec{x}^*) = 0$ and $\nabla^2 f(\vec{x}^*)$ is **positive definite**, then \vec{x}^* is a local minimizer.

Definition

- A matrix H is called **positive definite** if for any nonzero vector $\vec{v} \in \mathbb{R}^n$, $\vec{v}^\top H \vec{v} > 0$.
- H is called **positive semidefinite** if $\vec{v}^\top H \vec{v} \geq 0$ for all $\vec{v} \in \mathbb{R}^n$.
- H is **negative definite** or **negative semidefinite** if $-H$ is positive definite or positive semidefinite.
- H is **indefinite** if it is neither positive semidefinite nor negative semidefinite.

Convergence of the steepest descent method

Theorem (Convergence theorem of the steepest descent method)

If the steepest descent method converges to a local minimizer \vec{x}^* , where $\nabla^2 f(\vec{x})$ is positive definite, and e_{\max} and e_{\min} are the largest and the smallest eigenvalue of $\nabla^2 f(\vec{x})$, then

$$\lim_{k \rightarrow \infty} \frac{\|\vec{x}_{k+1} - \vec{x}^*\|}{\|\vec{x}_k - \vec{x}^*\|} \leq \left(\frac{e_{\max} - e_{\min}}{e_{\max} + e_{\min}} \right)$$

Definition

For a scalar λ and an unit vector v , (λ, v) is an eigenpair of of a matrix H if $Hv = \lambda v$. The scalar λ is called an eigenvalue of H , and v is called an eigenvector.

Newton's method

- We use the quadratic model to find the step length α_k . Can we use the quadratic model to find the search direction \vec{p}_k ?
- Yes, we can. Recall the quadratic model (now \vec{p} is the variable.)

$$f(\vec{x}_k + \vec{p}) \approx \frac{1}{2} \vec{p}^T H_k \vec{p} + \vec{p}^T \vec{g}_k + f_k$$

- Compute the gradient $\nabla_{\vec{p}} f(\vec{x}_k + \vec{p}) = H_k \vec{p} + \vec{g}_k$
- The solution of $\nabla_{\vec{p}} f(\vec{x}_k + \vec{p}) = 0$ is $\vec{p}_k = -H_k^{-1} \vec{g}_k$.
- Newton's method uses p_k as the search direction

Newton's method

- 1 Given an initial guess \vec{x}_0
- 2 For $k = 0, 1, 2, \dots$ until converge

$$\vec{x}_{k+1} = \vec{x}_k - H_k^{-1} \vec{g}_k.$$

Descent direction

- The direction $p_k = -H_k^{-1}g_k$ is called Newton's direction
- Is p_k a descent direction? (what's the definition of descent directions?)
- We only need to check if $\vec{g}_k^T \vec{p}_k < 0$.

$$\vec{g}_k^T \vec{p}_k = -\vec{g}_k^T H_k^{-1} \vec{g}_k.$$

Thus, \vec{p}_k is a descent direction if H^{-1} is positive definite.

- For a symmetric matrix H , the following conditions are equivalent
- H is positive definite.
 - H^{-1} is positive definite.
 - All the eigenvalues of H are positive.

Some properties of eigenvalues/eigenvectors

- A symmetric matrix H , of order n has n real eigenvalues and n real and linearly independent (orthogonal) eigenvectors

$$Hv_1 = \lambda_1 v_1, \quad Hv_2 = \lambda_2 v_2, \quad \dots, \quad Hv_n = \lambda_n v_n$$

- Let $V = [v_1 \ v_2 \ \dots \ v_n]$, $\Lambda = \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_n \end{bmatrix}$, $HV = V\Lambda$.

- If $\lambda_1, \lambda_2, \dots, \lambda_n$ are nonzero, since $H = V\Lambda V^{-1}$,

$$H^{-1} = V\Lambda^{-1}V^{-1}, \quad \Lambda^{-1} = \begin{bmatrix} 1/\lambda_1 & & & \\ & 1/\lambda_2 & & \\ & & \ddots & \\ & & & 1/\lambda_n \end{bmatrix}$$

The eigenvalues of H^{-1} are $\frac{1}{\lambda_1}, \frac{1}{\lambda_2}, \dots, \frac{1}{\lambda_n}$.

How to solve $H\vec{p} = -\vec{g}$?

- For a symmetric positive definite matrix H , $H\vec{p} = -\vec{g}$ can be solved by Cholesky decomposition, which is similar to LU decomposition, but is only half computational cost of LU decomposition.

- Let $H = \begin{bmatrix} h_{11} & h_{12} & h_{13} \\ h_{21} & h_{22} & h_{23} \\ h_{31} & h_{32} & h_{33} \end{bmatrix}$, where $h_{12} = h_{21}$, $h_{13} = h_{31}$, $h_{23} = h_{32}$.

Cholesky decomposition makes $H = LL^T$, where L is a lower

triangular matrix, $L = \begin{bmatrix} l_{11} & & \\ l_{21} & l_{22} & \\ l_{31} & l_{32} & l_{33} \end{bmatrix}$

- Using Cholesky decomposition, $H\vec{p} = -\vec{g}$ can be solved by
 - 1 Compute $H = LL^T$
 - 2 $\vec{p} = -(L^T)^{-1}L^{-1}\vec{g}$
- In Matlab, use $p = -H \setminus g$. Don't use $inv(H)$.

The Cholesky decomposition

For $i = 1, 2, \dots, n$

$$l_{ii} = \sqrt{h_{ii}}$$

For $j = i + 1, i + 2, \dots, n$

$$l_{ji} = \frac{h_{ji}}{l_{ii}}$$

For $k = i + 1, i + 2, \dots, j$

$$h_{jk} = h_{jk} - l_{ji}l_{ki}$$

$$\begin{bmatrix} h_{11} & h_{12} & h_{13} \\ h_{21} & h_{22} & h_{23} \\ h_{31} & h_{32} & h_{33} \end{bmatrix} = LL^T = \begin{bmatrix} l_{11}^2 & l_{11}l_{21} & l_{11}l_{31} \\ l_{11}l_{21} & l_{21}^2 + l_{22}^2 & l_{21}l_{31} + l_{22}l_{32} \\ l_{11}l_{31} & l_{21}l_{31} + l_{22}l_{32} & l_{31}^2 + l_{32}^2 + l_{33}^2 \end{bmatrix}$$

$$l_{11} = \sqrt{h_{11}} \quad h_{22}^{(2)} = h_{22} - l_{21}l_{21}$$

$$l_{21} = h_{21}/l_{11} \quad h_{32}^{(2)} = h_{32} - l_{21}l_{31}$$

$$l_{31} = h_{31}/l_{11} \quad h_{33}^{(2)} = h_{33} - l_{31}l_{31}$$

$$l_{22} = \sqrt{h_{22}^{(2)}}$$

$$l_{32} = h_{32}^{(2)}/l_{22}$$

$$l_{33} = \sqrt{h_{33}^{(2)} - l_{32}l_{32}}$$

Convergence of Newton's method

Theorem

Suppose f is twice differentiable. $\nabla^2 f$ is continuous in a neighborhood of \vec{x}^ and $\nabla^2 f(\vec{x}^*)$ is positive definite, and if \vec{x}_0 is sufficiently close to \vec{x}^* , the sequence converges to \vec{x}^* quadratically.*

Three problems of Newton's method

- 1 H may not be positive definite \Rightarrow Modified Newton's method + Line search.
- 2 H is expensive to compute \Rightarrow Quasi-Newton.
- 3 H^{-1} is expensive to compute \Rightarrow Conjugate gradient.