



Mathematical Foundations of Data Analysis (MFDA) - II

Boqiang Huang

huang@math.uni-koeln.de

Institute of Mathematics, University of Cologne, Germany

2019.04.09-11



Chapter 1. Basics

❖ 1. Probability and Information Theory

1.1. Basic definitions and rules in probability theory

1.2. Basic definitions and rules in information theory

❖ 2. Numerical computation

2.1 Basic knowledge in numerical computation

2.2 Basic knowledge in optimizations

2.2.1 Gradient-based optimization

2.2.2 Constrained optimization

❖ 3. Application: statistical model based data denoising (in tutorial after lecture this Thursday)

Reference:

[1] I. Goodfellow, Y. Bengio, A. Courville, Deep learning, Chapter 3-4, MIT Press, 2016.

[2] A. Antoniou, W.-S. Lu, Practical optimization: algorithms and engineering applications, Springer, 2007.

[3] I. Cohen, Noise spectrum estimation in adverse environments: improved minima controlled recursive averaging, IEEE Trans. Acoust., Speech, Signal Processing, vol. 11, no. 5, pp. 466-475, 2003.



2. Numerical Computation

- ❖ If there is no analytical solution, approximate/estimate it via iteratively numerical process

2.1 Overflow and Underflow

- Underflow occurs when numbers near zero are rounded to zero
 rounding error, avoid division by zero, ...
- Overflow occurs when numbers with large magnitude are approximated as ∞ or $-\infty$

Example: the softmax function

$$\text{softmax}(\mathbf{x})_i = \frac{\exp(x_i)}{\sum_{j=1}^n \exp(x_j)}$$



2. Numerical Computation

- ❖ If there is no analytical solution, approximate/estimate it via iteratively numerical process

2.2 Poor Conditioning

how rapidly a function changes with respect to small changes in its inputs

rounding errors in the inputs can result in large changes in the output

Example: consider a function $f(\mathbf{x}) = \mathbf{A}^{-1}\mathbf{x}$

condition number of matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ defined as ratio of the largest and smallest eigenvalue

$$\max_{i,j} \left| \frac{\lambda_i}{\lambda_j} \right|$$

when it is large, matrix inversion is sensitive to the input

Remark: poorly conditioned matrices amplify pre-existing errors when we multiply by its inverse



2. Numerical Computation

- ❖ If there is no analytical solution, approximate/estimate it via iteratively numerical process

2.3 Gradient-Based Optimization

$$\mathbf{x}^* = \arg \min f(\mathbf{x}) \quad \mathbf{x} \in \mathbb{R}^n$$

minimizer

objective/cost/loss/error function



2. Numerical Computation

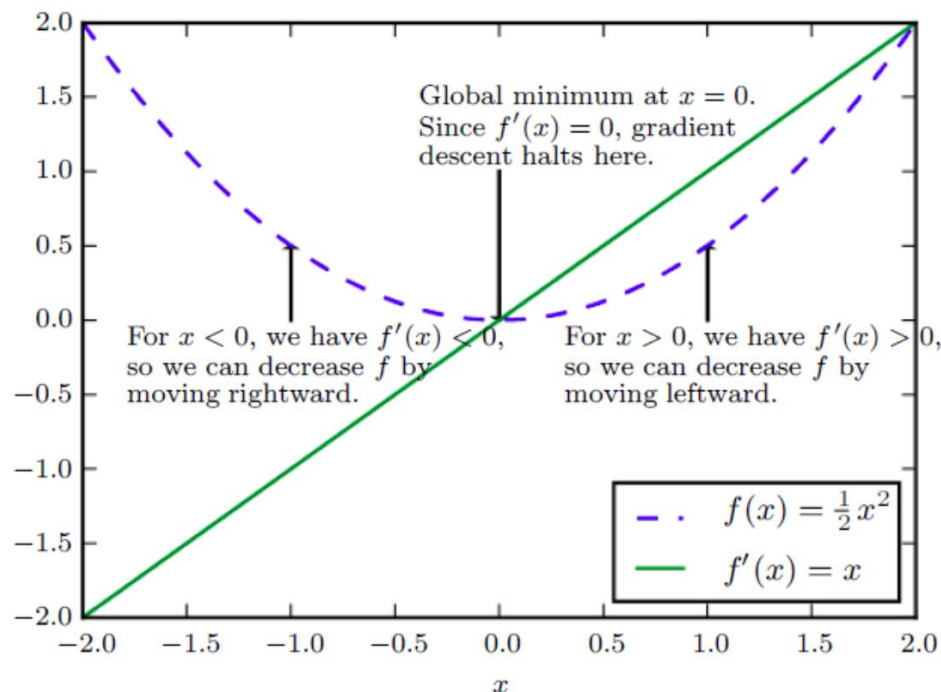
❖ If there is no analytical solution, approximate/estimate it via iteratively numerical process

2.3 Gradient-Based Optimization

$$\mathbf{x}^* = \arg \min f(\mathbf{x}) \quad \mathbf{x} \in \mathbb{R}^n$$

minimizer

objective/cost/loss/error function



Gradient Descent Method

consider a real-valued function $y = f(x)$

its derivative gives the slope of $f(x)$ at point x

$$f(x + \epsilon) \approx f(x) + \epsilon f'(x)$$

$f(x - \epsilon \text{sign}(f'(x)))$ is less than $f(x)$ for small enough ϵ

we can reduce $f(x)$ by moving x in small steps with opposite sign of the derivative



2. Numerical Computation

- ❖ If there is no analytical solution, approximate/estimate it via iteratively numerical process

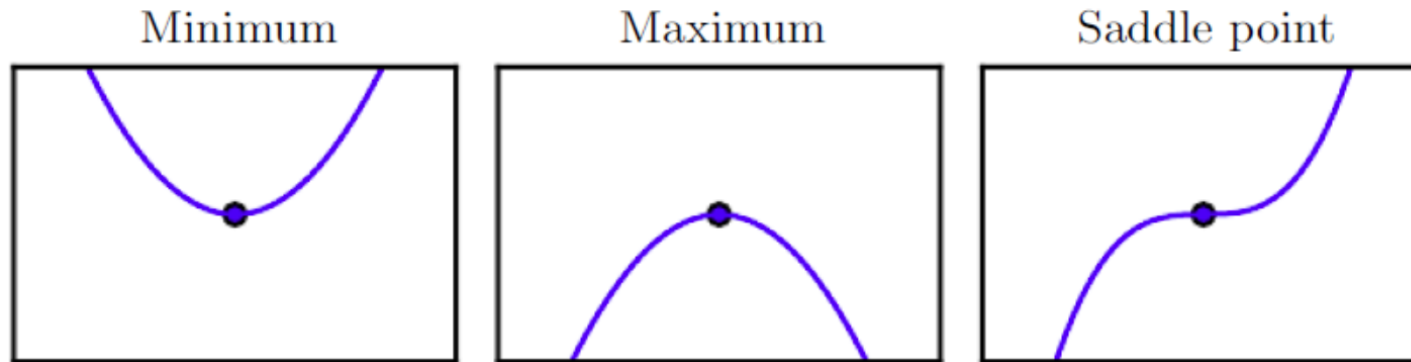
2.3 Gradient-Based Optimization

$$\mathbf{x}^* = \arg \min f(\mathbf{x}) \quad \mathbf{x} \in \mathbb{R}^n$$

minimizer

objective/cost/loss/error function

critical/stationary point exists at $f'(x) = 0$





2. Numerical Computation

- ❖ If there is no analytical solution, approximate/estimate it via iteratively numerical process

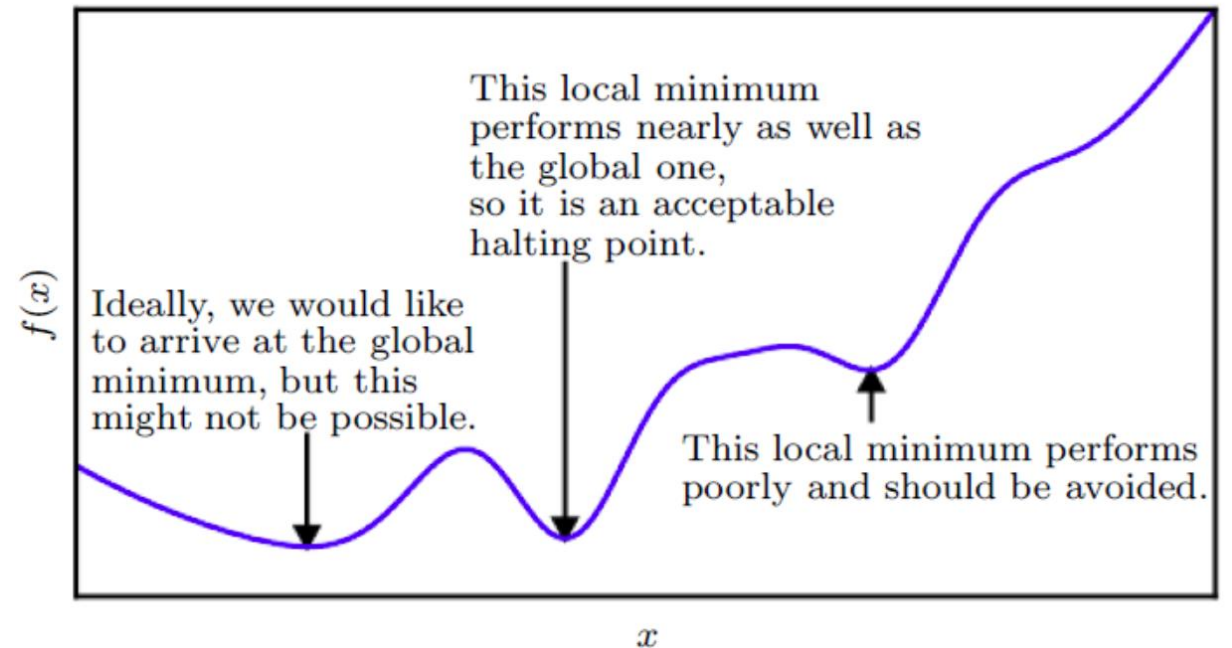
2.3 Gradient-Based Optimization

$$\mathbf{x}^* = \arg \min f(\mathbf{x}) \quad \mathbf{x} \in \mathbb{R}^n$$

minimizer

objective/cost/loss/error function

critical/stationary point exists at $f'(x) = 0$





2. Numerical Computation

- ❖ If there is no analytical solution, approximate/estimate it via iteratively numerical process

2.3 Gradient-Based Optimization

$$\mathbf{x}^* = \arg \min f(\mathbf{x}) \quad \mathbf{x} \in \mathbb{R}^n$$

minimizer

objective/cost/loss/error function

critical/stationary point exists at $\nabla_{\mathbf{x}} f(\mathbf{x}) = \left[\frac{\partial f}{\partial x_1} \quad \dots \quad \frac{\partial f}{\partial x_n} \right]^T = 0$

gradient

directional derivative in direction \mathbf{u} (a unit vector) is the slope of the function f in direction \mathbf{u} , $\mathbf{u} \in \mathbb{R}^n$

$$\nabla_{\mathbf{u}} f(\mathbf{x}) = \lim_{\alpha \rightarrow 0} \frac{f(\mathbf{x} + \alpha \mathbf{u}) - f(\mathbf{x})}{\alpha} \quad \left. \frac{\partial f(\mathbf{x} + \alpha \mathbf{u})}{\partial \alpha} \right|_{\alpha=0} = \mathbf{u}^T \nabla_{\mathbf{x}} f(\mathbf{x})$$



2. Numerical Computation

- ❖ If there is no analytical solution, approximate/estimate it via iteratively numerical process

2.3 Gradient-Based Optimization

$$\mathbf{x}^* = \arg \min f(\mathbf{x}) \quad \mathbf{x} \in \mathbb{R}^n$$

minimizer

objective/cost/loss/error function

critical/stationary point exists at $\nabla_{\mathbf{x}} f(\mathbf{x}) = \left[\frac{\partial f}{\partial x_1} \quad \dots \quad \frac{\partial f}{\partial x_n} \right]^T = 0$

gradient

directional derivative in direction \mathbf{u} (a unit vector) is the slope of the function f in direction \mathbf{u} , $\mathbf{u} \in \mathbb{R}^n$

$$\nabla_{\mathbf{u}} f(\mathbf{x}) = \lim_{\alpha \rightarrow 0} \frac{f(\mathbf{x} + \alpha \mathbf{u}) - f(\mathbf{x})}{\alpha}$$

$$\left. \frac{\partial f(\mathbf{x} + \alpha \mathbf{u})}{\partial \alpha} \right|_{\alpha=0} = \mathbf{u}^T \nabla_{\mathbf{x}} f(\mathbf{x})$$

$$\min_{\mathbf{u}, \mathbf{u}^T \mathbf{u} = 1} \mathbf{u}^T \nabla_{\mathbf{x}} f(\mathbf{x}) = \min_{\mathbf{u}, \mathbf{u}^T \mathbf{u} = 1} \|\mathbf{u}\|_2 \|\nabla_{\mathbf{x}} f(\mathbf{x})\|_2 \cos \theta$$

$$\min_{\mathbf{u}, \mathbf{u}^T \mathbf{u} = 1} \cos \theta$$

$$\mathbf{u} = -\nabla_{\mathbf{x}} f(\mathbf{x}) \quad \text{steepest/gradient descent}$$



2. Numerical Computation

- ❖ If there is no analytical solution, approximate/estimate it via iteratively numerical process

2.3 Gradient-Based Optimization

$$\mathbf{x}^* = \arg \min f(\mathbf{x}) \quad \mathbf{x} \in \mathbb{R}^n$$

minimizer

objective/cost/loss/error function

critical/stationary point exists at $\nabla_{\mathbf{x}} f(\mathbf{x}) = \left[\frac{\partial f}{\partial x_1} \quad \dots \quad \frac{\partial f}{\partial x_n} \right]^T = 0$

gradient

directional derivative in direction \mathbf{u} (a unit vector) is the slope of the function f in direction \mathbf{u} , $\mathbf{u} \in \mathbb{R}^n$

$$\mathbf{u} = -\nabla_{\mathbf{x}} f(\mathbf{x}) \quad \text{steepest/gradient descent}$$

iteration: $\mathbf{x}_{k+1} = \mathbf{x}_k - \epsilon \nabla_{\mathbf{x}} f(\mathbf{x}_k)$

ϵ : learning rate



2. Numerical Computation

- ❖ If there is no analytical solution, approximate/estimate it via iteratively numerical process

2.3 Gradient-Based Optimization

$$\mathbf{x}^* = \arg \min_{\mathbf{x}} f(\mathbf{x})$$

minimizer

objective/cost/loss/

critical/stationary point exists at $\nabla_{\mathbf{x}} f(\mathbf{x}) = \left[\frac{\partial f}{\partial x_1} \quad \dots \quad \frac{\partial f}{\partial x_n} \right]^T$

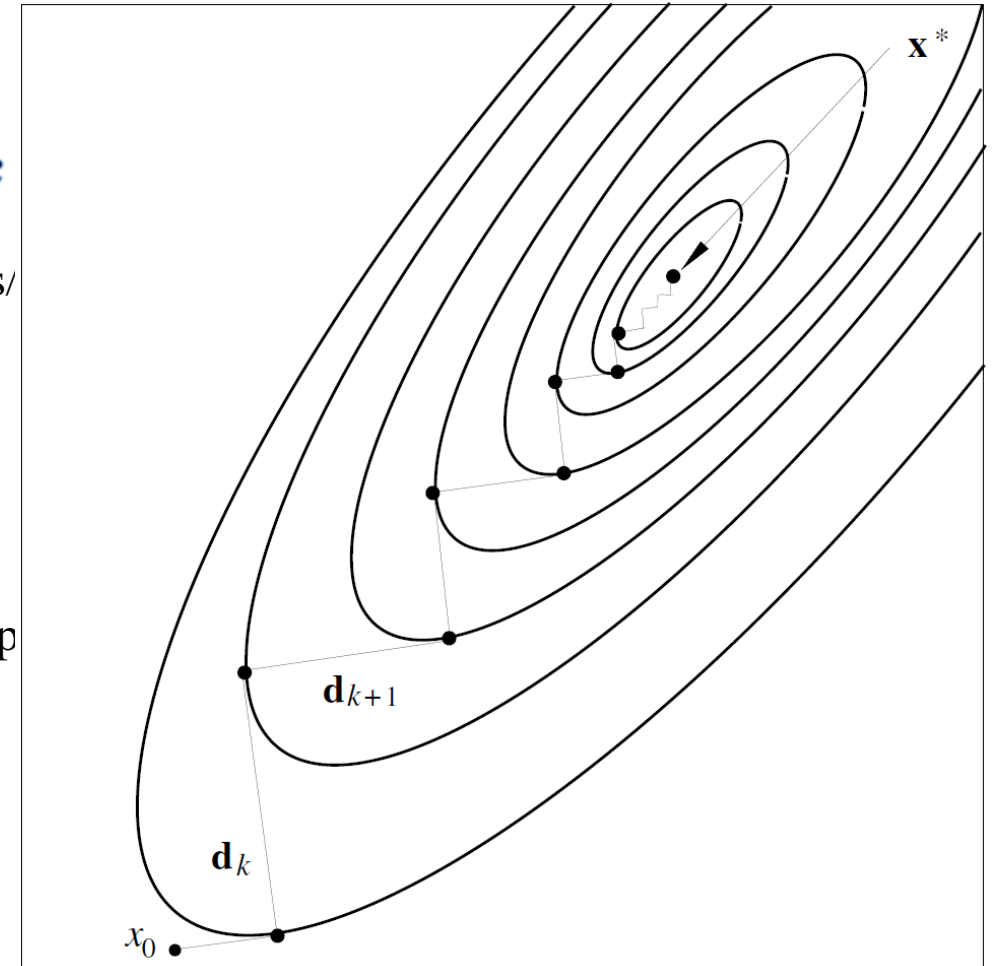
gradient

directional derivative in direction \mathbf{u} (a unit vector) is the slope

$$\mathbf{u} = -\nabla_{\mathbf{x}} f(\mathbf{x}) \quad \text{steepest/gradient descent}$$

iteration: $\mathbf{x}_{k+1} = \mathbf{x}_k - \epsilon \nabla_{\mathbf{x}} f(\mathbf{x}_k)$

ϵ : learning rate





2. Numerical Computation

- ❖ If there is no analytical solution, approximate/estimate it via iteratively numerical process

2.3 Gradient-Based Optimization

$$\mathbf{x}^* = \arg \min f(\mathbf{x}) \quad \mathbf{x} \in \mathbb{R}^n$$

minimizer

objective/cost/loss/error function

critical/stationary point exists at $\nabla_{\mathbf{x}} f(\mathbf{x}) = \left[\frac{\partial f}{\partial x_1} \quad \dots \quad \frac{\partial f}{\partial x_n} \right]^T = 0$

gradient

directional derivative in direction \mathbf{u} (a unit vector) is the slope of the function f in direction \mathbf{u} , $\mathbf{u} \in \mathbb{R}^n$

$$\mathbf{u} = -\nabla_{\mathbf{x}} f(\mathbf{x}) \quad \text{steepest/gradient descent}$$

iteration: $\mathbf{x}_{k+1} = \mathbf{x}_k - \epsilon_k \nabla_{\mathbf{x}} f(\mathbf{x}_k)$

line search: $\epsilon_k = \operatorname{argmin} f(\mathbf{x}_k + \epsilon \nabla_{\mathbf{x}} f(\mathbf{x}_k))$

ϵ_k : learning rate



2. Numerical Computation

- ❖ If there is no analytical solution, approximate/estimate it via iteratively numerical process

2.3 Gradient-Based Optimization

2.3.1 Beyond the Gradient: Jacobian and Hessian Matrices

$f: \mathbb{R}^m \rightarrow \mathbb{R}^n$ Jacobian Matrix $\mathbf{J} \in \mathbb{R}^{n \times m}$, where $J_{i,j} = \frac{\partial}{\partial x_j} f_i(\mathbf{x})$

$f: \mathbb{R}^m \rightarrow \mathbb{R}$ Hessian Matrix $\mathbf{H} \in \mathbb{R}^{m \times m}$, where $H_{i,j} = \frac{\partial^2}{\partial x_i \partial x_j} f(\mathbf{x})$



2. Numerical Computation

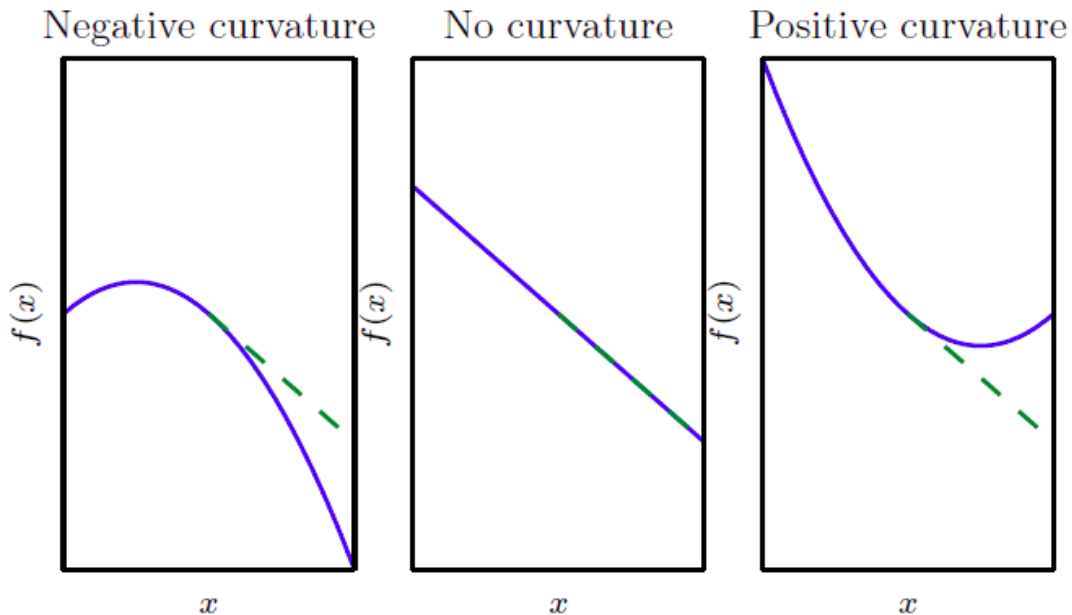
- ❖ If there is no analytical solution, approximate/estimate it via iteratively numerical process

2.3 Gradient-Based Optimization

2.3.1 Beyond the Gradient: Jacobian and Hessian Matrices

$$f: \mathbb{R}^m \rightarrow \mathbb{R}^n \quad \text{Jacobian Matrix } \mathbf{J} \in \mathbb{R}^{n \times m}, \text{ where } J_{i,j} = \frac{\partial}{\partial x_j} f_i(\mathbf{x})$$

$$f: \mathbb{R}^m \rightarrow \mathbb{R} \quad \text{Hessian Matrix } \mathbf{H} \in \mathbb{R}^{m \times m}, \text{ where } H_{i,j} = \frac{\partial^2}{\partial x_i \partial x_j} f(\mathbf{x})$$



second derivative, curvature

positive/negative (semi-) definite

function f is convex/concave



2. Numerical Computation

- ❖ If there is no analytical solution, approximate/estimate it via iteratively numerical process

2.3 Gradient-Based Optimization

2.3.1 Beyond the Gradient: Jacobian and Hessian Matrices

$f: \mathbb{R}^m \rightarrow \mathbb{R}^n$ Jacobian Matrix $\mathbf{J} \in \mathbb{R}^{n \times m}$, where $J_{i,j} = \frac{\partial}{\partial x_j} f_i(\mathbf{x})$

$f: \mathbb{R}^m \rightarrow \mathbb{R}$ Hessian Matrix $\mathbf{H} \in \mathbb{R}^{m \times m}$, where $H_{i,j} = \frac{\partial^2}{\partial x_i \partial x_j} f(\mathbf{x})$

consider a second-order Taylor approximation of $f(\mathbf{x})$

$$f(\mathbf{x}) \approx f(\mathbf{x}^{(0)}) + (\mathbf{x} - \mathbf{x}^{(0)})^\top \mathbf{g} + \frac{1}{2} (\mathbf{x} - \mathbf{x}^{(0)})^\top \mathbf{H} (\mathbf{x} - \mathbf{x}^{(0)})$$

where \mathbf{g} is the gradient and \mathbf{H} is the Hessian at $\mathbf{x}^{(0)}$. If we set $\mathbf{x} = \mathbf{x}^{(0)} - \epsilon \mathbf{g}$

$$f(\mathbf{x}^{(0)} - \epsilon \mathbf{g}) \approx f(\mathbf{x}^{(0)}) - \epsilon \mathbf{g}^\top \mathbf{g} + \frac{1}{2} \epsilon^2 \mathbf{g}^\top \mathbf{H} \mathbf{g}$$

if $\mathbf{g}^\top \mathbf{H} \mathbf{g}$ is positive

$$\epsilon^* = \frac{\mathbf{g}^\top \mathbf{g}}{\mathbf{g}^\top \mathbf{H} \mathbf{g}}$$



2. Numerical Computation

- ❖ If there is no analytical solution, approximate/estimate it via iteratively numerical process

2.3 Gradient-Based Optimization

2.3.1 Beyond the Gradient: Jacobian and Hessian Matrices

$$f: \mathbb{R}^m \rightarrow \mathbb{R}^n \quad \text{Jacobian Matrix } \mathbf{J} \in \mathbb{R}^{n \times m}, \text{ where } J_{i,j} = \frac{\partial}{\partial x_j} f_i(\mathbf{x})$$

$$f: \mathbb{R}^m \rightarrow \mathbb{R} \quad \text{Hessian Matrix } \mathbf{H} \in \mathbb{R}^{m \times m}, \text{ where } H_{i,j} = \frac{\partial^2}{\partial x_i \partial x_j} f(\mathbf{x})$$

consider a second-order Taylor approximation of $f(\mathbf{x})$

$$f(\mathbf{x}) \approx f(\mathbf{x}^{(0)}) + (\mathbf{x} - \mathbf{x}^{(0)})^\top \mathbf{g} + \frac{1}{2} (\mathbf{x} - \mathbf{x}^{(0)})^\top \mathbf{H} (\mathbf{x} - \mathbf{x}^{(0)})$$

where \mathbf{g} is the gradient and \mathbf{H} is the Hessian at $\mathbf{x}^{(0)}$. If we set $\mathbf{x} = \mathbf{x}^{(0)} - \epsilon \mathbf{g}$

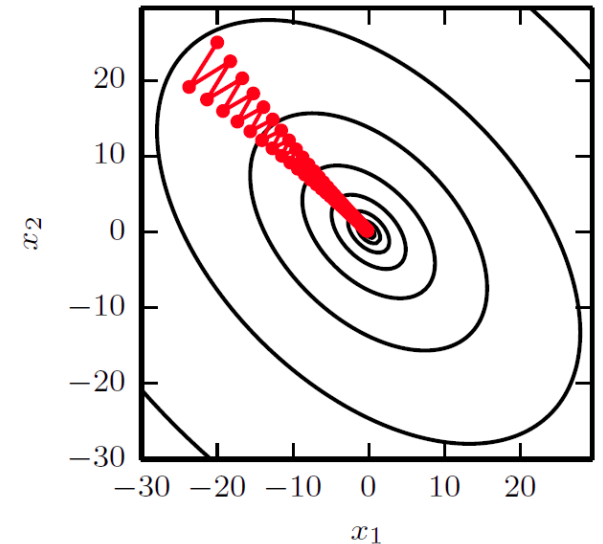
$$f(\mathbf{x}^{(0)} - \epsilon \mathbf{g}) \approx f(\mathbf{x}^{(0)}) - \epsilon \mathbf{g}^\top \mathbf{g} + \frac{1}{2} \epsilon^2 \mathbf{g}^\top \mathbf{H} \mathbf{g}$$

if $\mathbf{g}^\top \mathbf{H} \mathbf{g}$ is positive

$$\epsilon^* = \frac{\mathbf{g}^\top \mathbf{g}}{\mathbf{g}^\top \mathbf{H} \mathbf{g}}$$

poor conditioned Hessian leads to a poor gradient descent!!!

2019.04.09-11





2. Numerical Computation

- ❖ If there is no analytical solution, approximate/estimate it via iteratively numerical process

2.3 Gradient-Based Optimization

2.3.1 Beyond the Gradient: Jacobian and Hessian Matrices

Newton's method

$$f(\mathbf{x}) \approx f(\mathbf{x}^{(0)}) + (\mathbf{x} - \mathbf{x}^{(0)})^\top \nabla_{\mathbf{x}} f(\mathbf{x}^{(0)}) + \frac{1}{2} (\mathbf{x} - \mathbf{x}^{(0)})^\top \mathbf{H}(f)(\mathbf{x}^{(0)}) (\mathbf{x} - \mathbf{x}^{(0)})$$

solve for the critical point of this function

$$\mathbf{x}^* = \mathbf{x}^{(0)} - \mathbf{H}(f)(\mathbf{x}^{(0)})^{-1} \nabla_{\mathbf{x}} f(\mathbf{x}^{(0)})$$

Extension:

Gauss-Newton method

Conjugate-direction method

Quasi-Newton method

[2] A. Antoniou, W.-S. Lu, Practical optimization: algorithms and engineering applications, Springer, 2007.



2. Numerical Computation

- ❖ If there is no analytical solution, approximate/estimate it via iteratively numerical process

2.4 Constrained Optimization

minimize $f(\mathbf{x})$

subject to: $a_i(\mathbf{x}) = 0$ for $i = 1, 2, \dots, p$
 $c_j(\mathbf{x}) \geq 0$ for $j = 1, 2, \dots, q$

feasible region $\mathcal{R} = \{\mathbf{x} : a_i(\mathbf{x}) = 0 \text{ for } i = 1, 2, \dots, p, c_j(\mathbf{x}) \geq 0 \text{ for } j = 1, 2, \dots, q\}$

Karush-Kuhn-Tucker (KKT) approach, or generalized Lagrangian approach

$$L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\alpha}) = f(\mathbf{x}) + \sum_i \lambda_i a_i(\mathbf{x}) + \sum_j \alpha_j c_j(\mathbf{x})$$

First-order sufficient/necessary conditions for a minimum

Second-order sufficient/necessary conditions for a minimum

[2] A. Antoniou, W.-S. Lu, Practical optimization: algorithms and engineering applications, Springer, 2007.