

矩阵论

2024年秋季学期

第八讲

2024年10月9日

第4章 梯度分析与最优化

单变量函数

单变量函数的平稳点与极值点

平稳点: $f'(c) = 0$

- 平稳点是函数图像上的一个点, 在该点处, 函数的导数为零。直观上, 这意味着函数在这一点处的切线是水平的。
- 平稳点是函数局部极值可能出现的地方, 但不是所有平稳点都是极值点。

局部极小点:

$$f'(c) = 0 \quad f''(c) = \left. \frac{d^2 f(x)}{dx^2} \right|_{x=c} \geq 0$$

局部极大点:

$$f'(c) = 0 \quad f''(c) = \left. \frac{d^2 f(x)}{dx^2} \right|_{x=c} \leq 0$$

鞍点 (saddle point) : $f'(c) = 0$

$$f''(c + \Delta x) \leq 0 \quad f''(c + \Delta x) \geq 0$$

多变量函数

多变量函数的平稳点与极值点

多变量函数无约束极小化问题

$$\min_{\mathbf{x} \in S} f(\mathbf{x}) \quad f(\mathbf{x}): \mathbb{R}^n \rightarrow \mathbb{R}$$

开邻域

$$B(\mathbf{c}; r) = \{ \mathbf{x} \mid \mathbf{x} \in \mathbb{R}^n, \|\mathbf{x} - \mathbf{c}\|_2 < r \}$$

闭合邻域

$$B(\mathbf{c}; r) = \{ \mathbf{x} \mid \mathbf{x} \in \mathbb{R}^n, \|\mathbf{x} - \mathbf{c}\|_2 \leq r \}$$

多变量函数

多变量函数的平稳点与极值点

二阶泰勒级数逼近

$$\begin{aligned} f(\mathbf{c} + \Delta \mathbf{x}) &= f(\mathbf{c}) + \left(\frac{\partial f(\mathbf{c})}{\partial \mathbf{c}} \right)^T \Delta \mathbf{x} + \frac{1}{2} (\Delta \mathbf{x})^T \cdot \frac{\partial^2 f(\mathbf{c})}{\partial \mathbf{c} \partial \mathbf{c}^T} \Delta \mathbf{x} \\ &= f(\mathbf{c}) + (\nabla f(\mathbf{c}))^T \Delta \mathbf{x} + \frac{1}{2} (\Delta \mathbf{x})^T \mathbf{H}(f(\mathbf{c})) \Delta \mathbf{x} \end{aligned}$$

梯度向量

$$\nabla f(\mathbf{c}) = \frac{\partial f(\mathbf{c})}{\partial \mathbf{c}} = \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} \Big|_{\mathbf{x}=\mathbf{c}}$$

Hessian矩阵

$$\mathbf{H}(f(\mathbf{c})) = \frac{\partial^2 f(\mathbf{c})}{\partial \mathbf{c} \partial \mathbf{c}^T} = \frac{\partial^2 f(\mathbf{x})}{\partial \mathbf{x} \partial \mathbf{x}^T} \Big|_{\mathbf{x}=\mathbf{c}}$$

多变量函数

多变量函数的平稳点与极值点

$$f(c) \leq f(c + \Delta x) \quad \forall 0 < \|\Delta x\|_2 \leq r$$

局部极小

$$H(f(c)) = \frac{\partial^2 f(x)}{\partial x \partial x^T} \bigg|_{x=c} \succeq 0$$

局部极小

$$f(c) < f(c + \Delta x) \quad \forall 0 < \|\Delta x\|_2 \leq r$$

严格局部极小

$$f(c) \leq f(x) \quad \forall x \in S$$

全局极小

$$f(c) < f(x) \quad \forall x \in S, x \neq c$$

严格全局极小

多变量函数

多变量函数 $f(\mathbf{X})$ 的平稳点与极值点

邻域

$$B(\mathbf{C}; r) = \left\{ \mathbf{X} \mid \mathbf{X} \in \mathbb{R}^{m \times n}, \|\text{vec}(\mathbf{X}) - \text{vec}(\mathbf{C})\|_2 < r \right\}$$

二阶泰勒级数逼近

$$\begin{aligned} f(\mathbf{C} + \Delta \mathbf{X}) &= f(\mathbf{C}) + \left(\frac{\partial f(\mathbf{C})}{\partial \text{vec}(\mathbf{C})} \right)^T \text{vec}(\Delta \mathbf{X}) + \frac{1}{2} (\text{vec}(\Delta \mathbf{X}))^T \frac{\partial^2 f(\mathbf{C})}{\partial \text{vec}(\mathbf{C}) \partial (\text{vec} \mathbf{C})^T} \text{vec}(\Delta \mathbf{X}) \\ &= f(\mathbf{C}) + (\nabla_{\text{vec} \mathbf{C}} f(\mathbf{C}))^T \text{vec}(\Delta \mathbf{X}) + \frac{1}{2} (\text{vec}(\Delta \mathbf{X}))^T \mathbf{H}(f(\mathbf{C})) \text{vec}(\Delta \mathbf{X}) \end{aligned}$$

$$\nabla_{\text{vec} \mathbf{C}} f(\mathbf{C}) = \left. \frac{\partial f(\mathbf{X})}{\partial \text{vec}(\mathbf{X})} \right|_{\mathbf{X}=\mathbf{C}} \in \mathbb{R}^{mn} \quad \mathbf{H}(f(\mathbf{C})) = \left. \frac{\partial^2 f(\mathbf{X})}{\partial \text{vec}(\mathbf{X}) \partial (\text{vec} \mathbf{X})^T} \right|_{\mathbf{X}=\mathbf{C}} \in \mathbb{R}^{mn \times mn}$$

总结

表 4.1.1 实变函数的平稳点和极值点的条件

实变函数	$f(x) : \mathbb{R} \rightarrow \mathbb{R}$	$f(\mathbf{x}) : \mathbb{R}^n \rightarrow \mathbb{R}$	$f(\mathbf{X}) : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$
平稳点	$\left. \frac{\partial f(x)}{\partial x} \right _{x=c} = 0$	$\left. \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} \right _{\mathbf{x}=\mathbf{c}} = 0$	$\left. \frac{\partial f(\mathbf{X})}{\partial \mathbf{X}} \right _{\mathbf{X}=\mathbf{C}} = \mathbf{O}_{m \times n}$
局部极小点	$\left. \frac{\partial^2 f(x)}{\partial x \partial x} \right _{x=c} \geq 0$	$\left. \frac{\partial^2 f(\mathbf{x})}{\partial \mathbf{x} \partial \mathbf{x}^T} \right _{\mathbf{x}=\mathbf{c}} \succeq 0$	$\left. \frac{\partial^2 f(\mathbf{X})}{\partial \text{vec}(\mathbf{X}) \partial (\text{vec} \mathbf{X})^T} \right _{\mathbf{X}=\mathbf{C}} \succeq 0$
严格局部极小点	$\left. \frac{\partial^2 f(x)}{\partial x \partial x} \right _{x=c} > 0$	$\left. \frac{\partial^2 f(\mathbf{x})}{\partial \mathbf{x} \partial \mathbf{x}^T} \right _{\mathbf{x}=\mathbf{c}} \succ 0$	$\left. \frac{\partial^2 f(\mathbf{X})}{\partial \text{vec}(\mathbf{X}) \partial (\text{vec} \mathbf{X})^T} \right _{\mathbf{X}=\mathbf{C}} \succ 0$
局部极大点	$\left. \frac{\partial^2 f(x)}{\partial x \partial x} \right _{x=c} \leq 0$	$\left. \frac{\partial^2 f(\mathbf{x})}{\partial \mathbf{x} \partial \mathbf{x}^T} \right _{\mathbf{x}=\mathbf{c}} \preceq 0$	$\left. \frac{\partial^2 f(\mathbf{X})}{\partial \text{vec}(\mathbf{X}) \partial (\text{vec} \mathbf{X})^T} \right _{\mathbf{X}=\mathbf{C}} \preceq 0$
严格局部极大点	$\left. \frac{\partial^2 f(x)}{\partial x \partial x} \right _{x=c} < 0$	$\left. \frac{\partial^2 f(\mathbf{x})}{\partial \mathbf{x} \partial \mathbf{x}^T} \right _{\mathbf{x}=\mathbf{c}} \prec 0$	$\left. \frac{\partial^2 f(\mathbf{X})}{\partial \text{vec}(\mathbf{X}) \partial (\text{vec} \mathbf{X})^T} \right _{\mathbf{X}=\mathbf{C}} \prec 0$
鞍点	$\left. \frac{\partial^2 f(x)}{\partial x \partial x} \right _{x=c}$ 不定	$\left. \frac{\partial^2 f(\mathbf{x})}{\partial \mathbf{x} \partial \mathbf{x}^T} \right _{\mathbf{x}=\mathbf{c}}$ 不定	$\left. \frac{\partial^2 f(\mathbf{X})}{\partial \text{vec}(\mathbf{X}) \partial (\text{vec} \mathbf{X})^T} \right _{\mathbf{X}=\mathbf{C}}$ 不定

4.1.4 一阶必要条件, 二阶必要条件, 二阶充分条件

梯度分析与最优化

复变函数的平稳点和极值点条件

一阶微分

$$df(z, z^*) = \frac{\partial f(z, z^*)}{\partial \mathbf{z}^T} d\mathbf{z} + \frac{\partial f(z, z^*)}{\partial \mathbf{z}^H} d\mathbf{z}^* = \begin{bmatrix} \frac{\partial f(z, z^*)}{\partial \mathbf{z}^T} & \frac{\partial f(z, z^*)}{\partial \mathbf{z}^H} \end{bmatrix} \begin{bmatrix} d\mathbf{z} \\ d\mathbf{z}^* \end{bmatrix}$$

二阶微分

$$d^2 f(z, z^*) = \left(\frac{\partial^2 f(z, z^*)}{\partial \mathbf{z} \partial \mathbf{z}^T} d\mathbf{z} + \frac{\partial^2 f(z, z^*)}{\partial \mathbf{z}^* \partial \mathbf{z}^T} d\mathbf{z}^* \right)^T d\mathbf{z} + \left(\frac{\partial^2 f(z, z^*)}{\partial \mathbf{z} \partial \mathbf{z}^H} d\mathbf{z} + \frac{\partial^2 f(z, z^*)}{\partial \mathbf{z}^* \partial \mathbf{z}^H} d\mathbf{z}^* \right)^T d\mathbf{z}^*$$
$$= \begin{bmatrix} d\mathbf{z}^H & d\mathbf{z}^T \end{bmatrix} \begin{bmatrix} \frac{\partial^2 f(z, z^*)}{\partial \mathbf{z}^* \partial \mathbf{z}^T} & \frac{\partial^2 f(z, z^*)}{\partial \mathbf{z}^* \partial \mathbf{z}^H} \\ \frac{\partial^2 f(z, z^*)}{\partial \mathbf{z} \partial \mathbf{z}^T} & \frac{\partial^2 f(z, z^*)}{\partial \mathbf{z} \partial \mathbf{z}^H} \end{bmatrix} \begin{bmatrix} d\mathbf{z} \\ d\mathbf{z}^* \end{bmatrix}$$

梯度分析与最优化

复变函数的平稳点和极值点条件

二阶Taylor级数逼近

$$\begin{aligned} f(\mathbf{z}, \mathbf{z}^*) &\approx f(\mathbf{c}, \mathbf{c}^*) + \left[\frac{\partial f(\mathbf{c}, \mathbf{c}^*)}{\partial \mathbf{c}}, \frac{\partial f(\mathbf{c}, \mathbf{c}^*)}{\partial \mathbf{c}^*} \right] \begin{bmatrix} \Delta \mathbf{c} \\ \Delta \mathbf{c}^* \end{bmatrix} + \frac{1}{2} \begin{bmatrix} \Delta \mathbf{c}^H, \Delta \mathbf{c}^T \end{bmatrix} \begin{bmatrix} \frac{\partial^2 f(\mathbf{c}, \mathbf{c}^*)}{\partial \mathbf{c}^* \partial \mathbf{c}^T} & \frac{\partial^2 f(\mathbf{c}, \mathbf{c}^*)}{\partial \mathbf{c}^* \partial \mathbf{c}^H} \\ \frac{\partial^2 f(\mathbf{c}, \mathbf{c}^*)}{\partial \mathbf{c} \partial \mathbf{c}^T} & \frac{\partial^2 f(\mathbf{c}, \mathbf{c}^*)}{\partial \mathbf{c} \partial \mathbf{c}^H} \end{bmatrix} \begin{bmatrix} \Delta \mathbf{c} \\ \Delta \mathbf{c}^* \end{bmatrix} \\ &= f(\mathbf{c}, \mathbf{c}^*) + \left(\nabla f(\mathbf{c}, \mathbf{c}^*) \right)^T \Delta \tilde{\mathbf{c}} + \frac{1}{2} (\Delta \tilde{\mathbf{c}})^H \mathbf{H}(f(\mathbf{c}, \mathbf{c}^*)) \Delta \tilde{\mathbf{c}} \end{aligned}$$

梯度分析与最优化

极值点的辨识

必要条件

$$\left. \frac{\partial f(z, z^*)}{\partial z^*} \right|_{z=z_0} = 0$$

$$\begin{bmatrix} H_{z^*, z} & H_{z^*, z^*} \\ H_{z, z} & H_{z, z^*} \end{bmatrix}_{z=z_0} \succcurlyeq 0$$

或

$$\left. \frac{\partial f(Z, Z^*)}{\partial Z^*} \right|_{Z=Z_0} = 0$$

$$\begin{bmatrix} H_{Z^*, Z} & H_{Z^*, Z^*} \\ H_{Z, Z} & H_{Z, Z^*} \end{bmatrix}_{Z=Z_0} \succcurlyeq 0$$

充分条件

$$\left. \frac{\partial f(z, z^*)}{\partial z^*} \right|_{z=z_0} = 0$$

$$\begin{bmatrix} H_{z^*, z} & H_{z^*, z^*} \\ H_{z, z} & H_{z, z^*} \end{bmatrix}_{z=z_0} \succ 0$$

或

$$\left. \frac{\partial f(Z, Z^*)}{\partial Z^*} \right|_{Z=Z_0} = 0$$

$$\begin{bmatrix} H_{Z^*, Z} & H_{Z^*, Z^*} \\ H_{Z, Z} & H_{Z, Z^*} \end{bmatrix}_{Z=Z_0} \succ 0$$

则 z_0 为严格局部最小点

无约束最小化问题的梯度分析

给定一个实值目标函数 $f(w, w^*)$ 或 $f(W, W^*)$

无约束最小化问题的梯度分析

1. 共轭梯度矩阵决定最小化问题的闭式解。
2. 共轭梯度矩阵与 Hessian 矩阵给出局部极小点辨识的必要条件或充分条件。
3. 共轭梯度向量的负方向决定求解最小化问题的最速下降迭代算法。
4. Hessian 矩阵给出求解最小化问题的 Newton 算法。

无约束最小化问题的梯度分析——闭式解

例 4.2.1 考察求解超定矩阵方程 $Az = b$ 的最小二乘方法。定义误差平方和

$$\begin{aligned} J(z) &= \|Az - b\|_2^2 = (Az - b)^H (Az - b) \\ &= z^H A^H A z - z^H A^H b - b^H A z + b^H b \end{aligned}$$

为准则函数。令其共轭梯度向量 $\nabla_{z^*} J(z) = A^H A z - A^H b$ 等于零向量，易知：若 $A^H A$ 非奇异，则

$$z = (A^H A)^{-1} A^H b \quad \text{极大值/极小值?} \quad (4.2.24)$$

梯度向量

$$\nabla f(z, z^*) = \begin{bmatrix} \frac{\partial f(z, z^*)}{\partial z} \\ \frac{\partial f(z, z^*)}{\partial z^*} \end{bmatrix} \in \mathbb{C}^{2n}$$

无约束最小化问题的梯度分析——闭式解

例 4.2.2 考察求解超定矩阵方程 $Az = b$ 的最大似然方法。定义对数似然函数

$$l(\hat{z}) = C - \frac{1}{\sigma^2} e^H e = C - \frac{1}{\sigma^2} (b - A\hat{z})^H (b - A\hat{z}) \quad (4.2.25)$$

式中, C 为一实常数。求对数似然函数

$$l(\hat{z}) = C - \frac{1}{\sigma^2} b^H b + \frac{1}{\sigma^2} b^H A\hat{z} + \frac{1}{\sigma^2} \hat{z}^H A^H b - \frac{1}{\sigma^2} \hat{z}^H A^H A\hat{z} \quad (4.2.26)$$

相对于 z 的共轭梯度, 得

$$\nabla_{\hat{z}} l(\hat{z}) = \frac{1}{\sigma^2} A^H b - \frac{1}{\sigma^2} A^H A\hat{z}$$

令其等于零, 得 $A^H b - A^H A z_{\text{opt}} = 0$ 或 $A^H A z_{\text{opt}} = A^H b$, 其中 z_{opt} 是使对数似然函数 $l(\hat{z})$ 极大化的 \hat{z} 值。于是, 若 $A^H A$ 非奇异, 则

$$z_{\text{opt}} = (A^H A)^{-1} A^H b \quad (4.2.27)$$

与最小二乘解具有等价性

无约束最小化问题的梯度分析——极值点的辨识

必要条件

$$\left. \frac{\partial f(z, z^*)}{\partial z^*} \right|_{z=z_0} = 0, \quad \begin{bmatrix} H_{z^*,z} & H_{z^*,z^*} \\ H_{z,z} & H_{z,z^*} \end{bmatrix}_{z=z_0} \succcurlyeq 0$$

或

$$\left. \frac{\partial f(Z, Z^*)}{\partial Z^*} \right|_{Z=Z_0} = 0, \quad \begin{bmatrix} H_{Z^*,Z} & H_{Z^*,Z^*} \\ H_{Z,Z} & H_{Z,Z^*} \end{bmatrix}_{Z=Z_0} \succcurlyeq 0$$

充分条件

$$\left. \frac{\partial f(z, z^*)}{\partial z^*} \right|_{z=z_0} = 0, \quad \begin{bmatrix} H_{z^*,z} & H_{z^*,z^*} \\ H_{z,z} & H_{z,z^*} \end{bmatrix}_{z=z_0} \succ 0$$

或

$$\left. \frac{\partial f(Z, Z^*)}{\partial Z^*} \right|_{Z=Z_0} = 0, \quad \begin{bmatrix} H_{Z^*,Z} & H_{Z^*,Z^*} \\ H_{Z,Z} & H_{Z,Z^*} \end{bmatrix}_{Z=Z_0} \succ 0$$

则 z_0 为严格局部最小点

无约束最小化问题的梯度分析——凸优化理论

标准的约束优化问题

$$\min_{\mathbf{x}} f_0(\mathbf{x}) \quad \text{subject to } f_i(\mathbf{x}) \leq 0, i = 1, \dots, m; \mathbf{Ax} = \mathbf{b}$$

约束优化问题很难求解，决策变量 \mathbf{x} 很大，有很多局部解，收敛速度差，收敛的停止准则失败等。

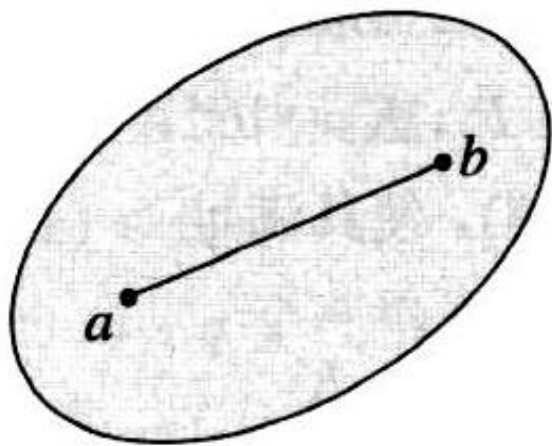
定义 4.3.1 一个集合 $S \in \mathbb{R}^n$ 称为凸集 (合)，若对任意两个点 $\mathbf{x}, \mathbf{y} \in S$ ，连接它们的线段也在集合 S 内，即

$$\mathbf{x}, \mathbf{y} \in S, \quad \theta \in [0, 1] \quad \Rightarrow \quad \theta \mathbf{x} + (1 - \theta) \mathbf{y} \in S \quad (4.3.11)$$

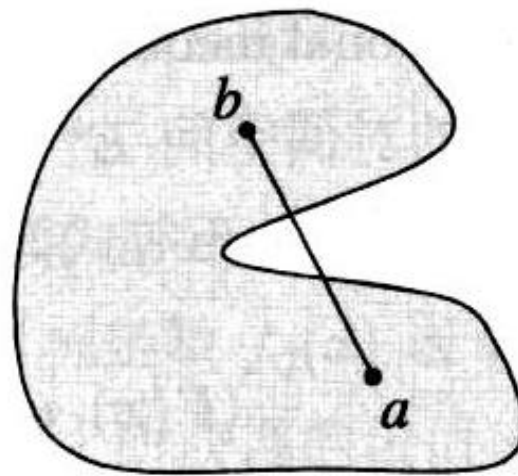
无约束最小化问题的梯度分析——凸优化理论

定义 4.3.1 一个集合 $S \in \mathbb{R}^n$ 称为凸集 (合), 若对任意两个点 $x, y \in S$, 连接它们的线段也在集合 S 内, 即

$$x, y \in S, \quad \theta \in [0, 1] \quad \Rightarrow \quad \theta x + (1 - \theta)y \in S \quad (4.3.11)$$



(a) 凸集



(b) 非凸集

无约束最小化问题的梯度分析——凸优化理论

凸函数

定义 4.3.4^[363] 给定一个凸集 $S \in \mathbb{R}^n$ 和函数 $f: S \rightarrow \mathbb{R}$, 则:

(1) 函数 $f: \mathbb{R}^n \rightarrow \mathbb{R}$ 称为凸函数 (convex function), 当且仅当 $S = \text{dom}(f)$ 是凸集, 并且对于所有 $\mathbf{x}, \mathbf{y} \in S$ 和每一个标量 $\alpha \in (0, 1)$, 函数满足 Jensen 不等式

$$f(\alpha \mathbf{x} + (1 - \alpha) \mathbf{y}) \leq \alpha f(\mathbf{x}) + (1 - \alpha) f(\mathbf{y}) \quad (4.3.18)$$

(2) 函数 $f(\mathbf{x})$ 称为严格凸函数 (strictly convex function), 当且仅当 $S = \text{dom}(f)$ 是凸集, 并且对于所有 $\mathbf{x}, \mathbf{y} \in S$ 和每一个标量 $\alpha \in (0, 1)$, 函数满足不等式

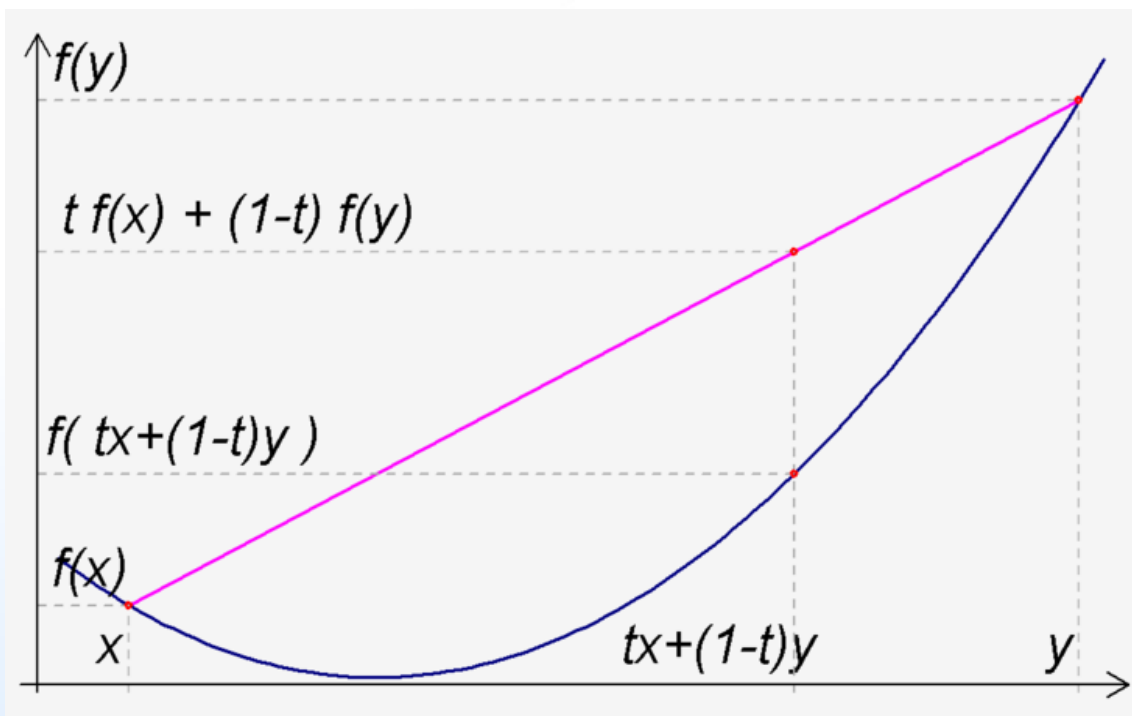
$$f(\alpha \mathbf{x} + (1 - \alpha) \mathbf{y}) < \alpha f(\mathbf{x}) + (1 - \alpha) f(\mathbf{y}) \quad (4.3.19)$$

无约束最小化问题的梯度分析—凸优化理论

凸函数

(2) 函数 $f(\mathbf{x})$ 称为严格凸函数 (strictly convex function), 当且仅当 $S = \text{dom}(f)$ 是凸集, 并且对于所有 $\mathbf{x}, \mathbf{y} \in S$ 和每一个标量 $\alpha \in (0, 1)$, 函数满足不等式

$$f(\alpha \mathbf{x} + (1 - \alpha) \mathbf{y}) < \alpha f(\mathbf{x}) + (1 - \alpha) f(\mathbf{y}) \quad (4.3.19)$$



无约束最小化问题的梯度分析—凸优化理论

凸函数辨识的二阶充分必要条件

定理 4.3.3^[328] 令 $f: S \rightarrow \mathbb{R}$ 是一个定义在 n 维向量空间 \mathbb{R}^n 内的凸集 S 上的函数, 并且可二次微分, 则 $f(x)$ 是凸函数, 当且仅当 Hessian 矩阵半正定

$$H_x f(x) = \frac{\partial^2 f(x)}{\partial x \partial x^T} \succeq 0, \quad \forall x \in S \quad (4.3.33)$$

注释 令 $f: S \rightarrow \mathbb{R}$ 是一个定义在 n 维向量空间 \mathbb{R}^n 内的凸集 S 上的函数, 并且可二次微分, 则 $f(x)$ 是严格凸函数, 当且仅当 Hessian 矩阵正定

$$H_x f(x) = \frac{\partial^2 f(x)}{\partial x \partial x^T} \succ 0, \quad \forall x \in S \quad (4.3.34)$$

与严格极小点的充分条件要求 Hessian 矩阵在 c 一点正定不同, 这里要求 Hessian 矩阵在整个凸集 S 的所有点均正定。

无约束最小化问题的梯度分析——实值目标函数的最速下降方向

以复向量或者矩阵为变元的实值目标函数的平稳点存在两种选择

$$\left. \frac{\partial f(Z, Z^*)}{\partial Z} \right|_{Z=C} = O_{m \times n} \quad \text{或} \quad \left. \frac{\partial f(Z, Z^*)}{\partial Z^*} \right|_{Z=C} = O_{m \times n}$$

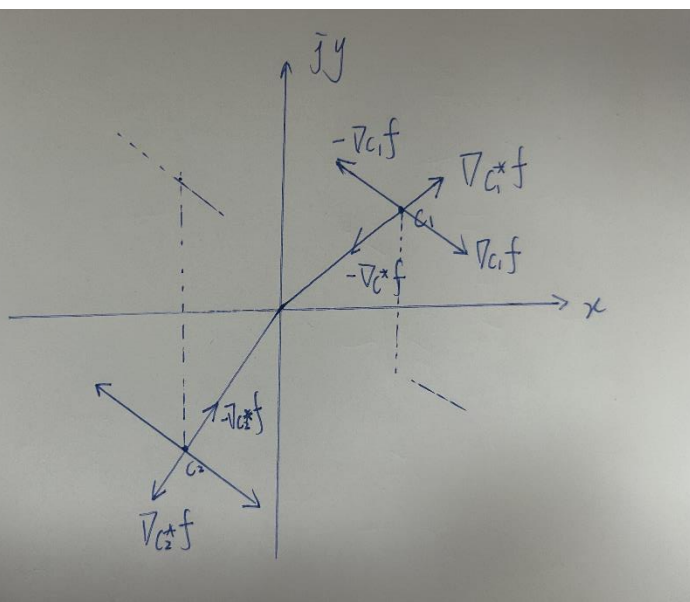
在设计优化迭代算法时，应该选哪一种梯度？

无约束最小化问题的梯度分析——实值目标函数的最速下降方向

定理 4.2.1^[58] 令 $f(z)$ 是复向量 z 的实值函数。通过将 z 和 z^* 视为独立的变元，实目标函数 $f(z)$ 的曲率方向由共轭梯度向量 $\nabla_{z^*} f(z)$ 给出。

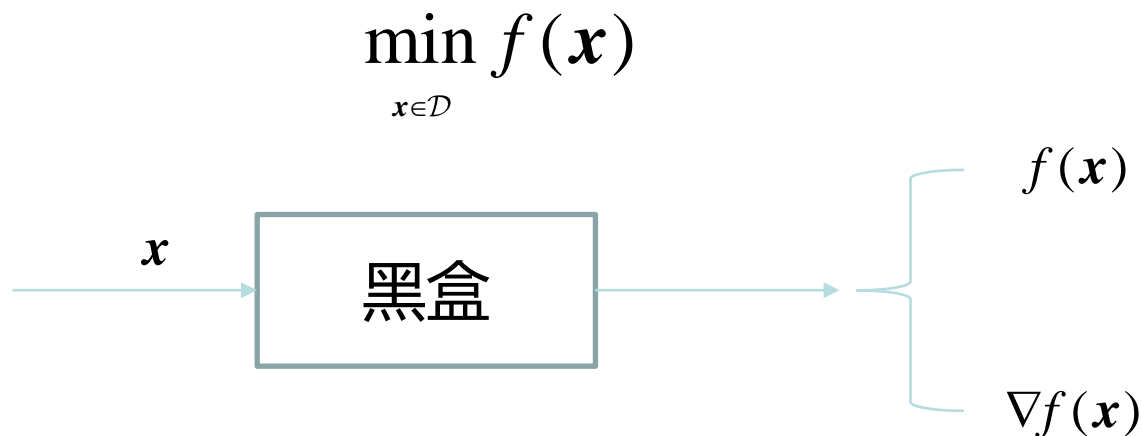
曲率方向也就是函数的最大变化率方向。

- (1) 共轭梯度向量 $\nabla_{z^*} f(z, z^*)$ 或 $\nabla_{\text{vec}(Z^*)} f(Z, Z^*)$ 给出目标函数增长最快的方向；
- (2) 负共轭梯度向量 $-\nabla_{z^*} f(z, z^*)$ 或 $-\nabla_{\text{vec}(Z^*)} f(Z, Z^*)$ 给出目标函数最陡减小的方向。



函数 $f(z) = |z|^2$ 梯度和共轭梯度方向示意图

平滑凸优化的一阶算法——梯度法/最陡下降方法



令 x_{opt} 表示 $\min f(x)$ 的最优解, 一阶黑盒优化 (first-order black-box optimization) 就是只利用 $f(x)$ 和 $\nabla f(x)$ 求解向量 $y \in \mathcal{Q}$ 满足 $y: f(y) - f(x_{\text{opt}}) \leq \varepsilon \longrightarrow$ 给定的精度误差

$$\arg \min f(x)$$

↓
Argument, 函数取最小值时 x 的取值

平滑凸优化的一阶算法——梯度法/最陡下降方法

梯度法（最陡下降法）是一种最简单的一阶优化算法

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \mu_k \Delta \mathbf{x}_k, \quad k = 1, 2, \dots$$



\mathbf{x}_{opt}



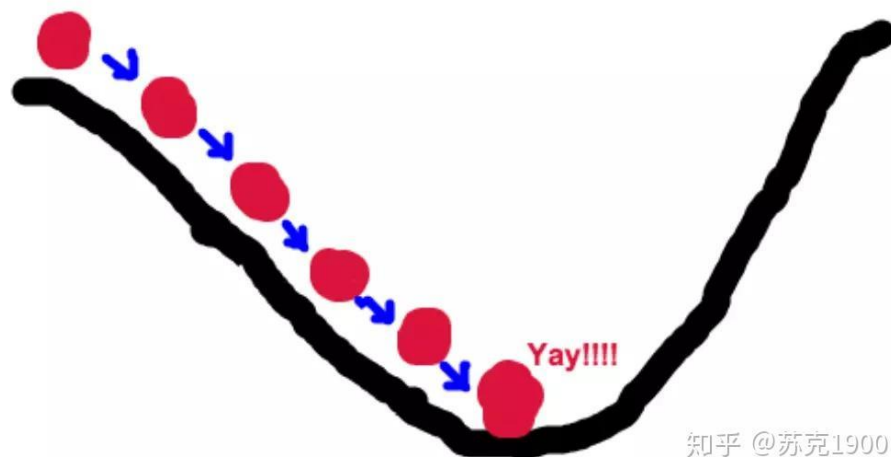
搜索方向或者更新方向

步长(学习率)
用于控制更新 \mathbf{x} 寻优的步
伐

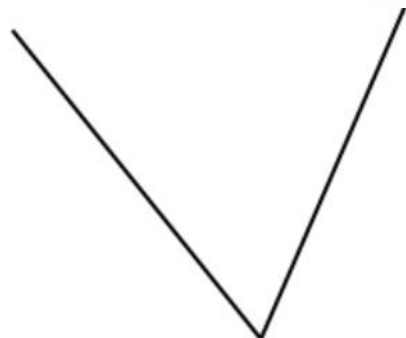
$$f(\mathbf{x}_{k+1}) < f(\mathbf{x}_k)$$

迭代过程中目标函数下降

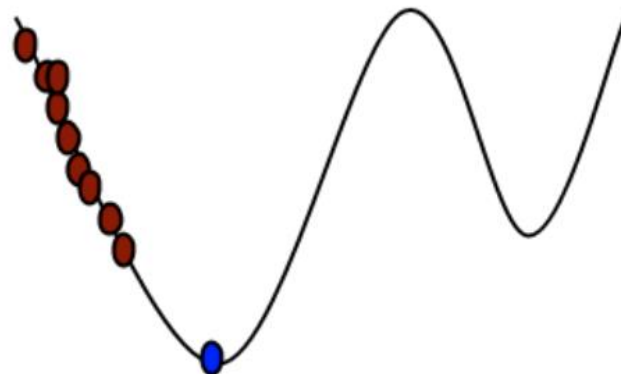
平滑凸优化的一阶算法——梯度法/最陡下降方法



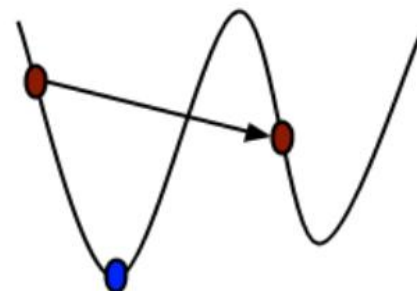
Gradient Descent
lecture notes from
UD262 Udacity Georgia
Tech ML Course.



not differentiable at
the corner



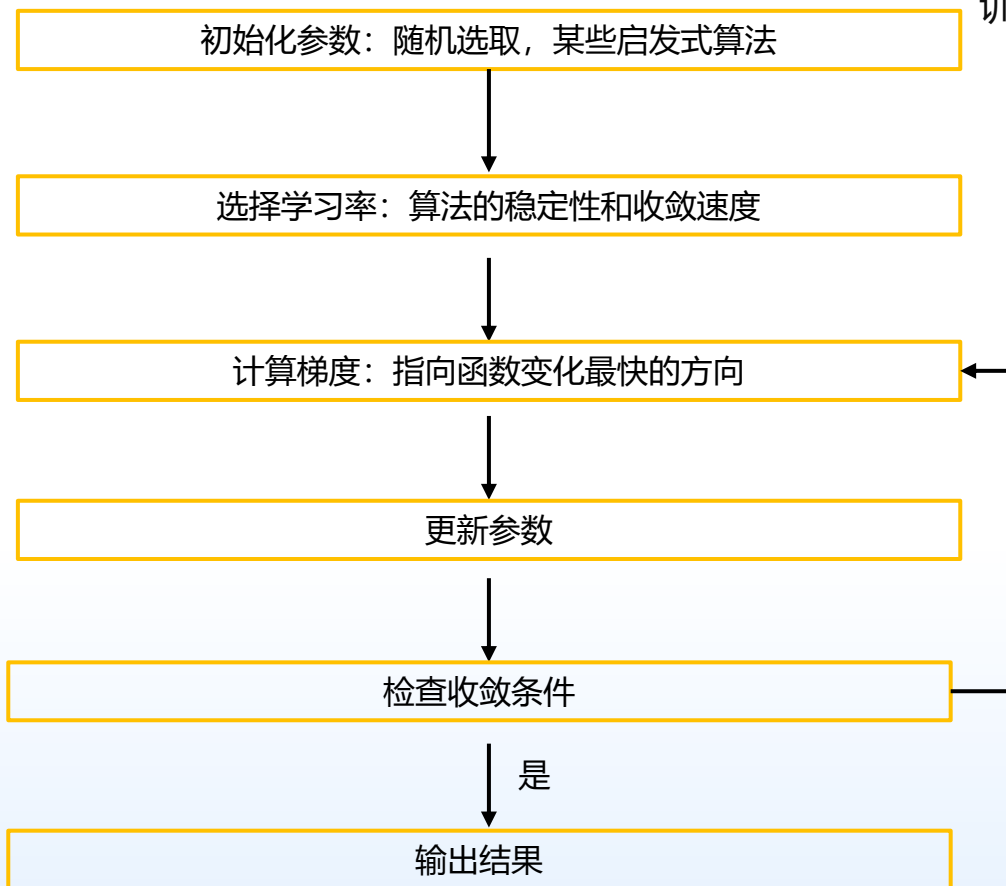
very small learning
rate needs lots of
steps



too big learning rate:
missed the minimum

平滑凸优化的一阶算法——梯度法/最陡下降方法

梯度法的基本流程



梯度消失和梯度爆炸：在深度学习模型训练中，由于层数过多，梯度可能在传播过程中逐渐变小或变大，导致训练难以进行。需要采取特定策略来解决这些问题。

拓展——梯度下降的变体

- 批量梯度下降 (Batch Gradient Descent)：每一步更新都使用所有训练样本来计算梯度。精度高但计算量大，对大数据集不够高效。
- 随机梯度下降 (Stochastic Gradient Descent, SGD)：每一步更新只使用一个训练样本来计算梯度。计算速度快，但更新过程中有较多噪声。
- 小批量梯度下降 (Mini-batch Gradient Descent)：每一步更新使用一小批训练样本来计算梯度。兼顾了批量梯度下降的精度和随机梯度下降的速度，是实际应用中的常用选择。

第三章习题

见学在浙大
作业版块
Homework2

10.13 交