

智能优化方法及其应用

第一讲 经典优化算法1

课程纲要

- 最优化理论数学基础复习
- 线搜索方法
- 梯度法和牛顿法
- 共轭梯度法
- 拟牛顿法
- 最小二乘法
- 最优性条件(*)
- 线性规划
- 二次规划

课程纲要

- 最优化理论数学基础复习
- 线搜索方法
- 梯度法和牛顿法
- 共轭梯度法
- 拟牛顿法
- 最小二乘法
- 最优性条件(*)
- 线性规划
- 二次规划

最优化理论数学基础

- 最优化问题就是求一个多元函数在某个给定集合上的极值问题, 几乎所有类型的最优化问题都可以描述如下:

$$\begin{aligned} \min f(x), \\ \text{s. t. } x \in \Omega \end{aligned}$$

其中: Ω 为某个给定的集合, 称为可行集或可行域;

$f(x)$ 为定义在集合 Ω 上连续可微的多元实值函数, 称为目标函数; $x = (x_1, x_2, \dots, x_n)^T$ 为决策变量; s.t. 为 **subject to** (受限于) 的缩写.

对于极大化问题, 在目标函数前添加负号可以转化为极小化问题. 因此, 只需考虑目标函数极小化问题.

最优化理论数学基础

- 最优化问题:

$$\min f(\mathbf{x}),$$

$$\text{s. t. } \mathbf{x} \in \Omega$$

可行域 Ω 一般常用等式和不等式来描述

$$\Omega = \{\mathbf{x} \in \mathbb{R}^n \mid c_i(\mathbf{x}) = 0, i \in \mathcal{E}; c_i(\mathbf{x}) \geq 0, i \in \mathcal{J}\}$$

其中 $c_i(\mathbf{x})$ ($i \in \mathcal{E} \cup \mathcal{J}$)为定义在 \mathbb{R}^n 上连续可微的多元实值函数,称为约束函数.

等式约束: $c_i(\mathbf{x}) = 0, i \in \mathcal{E}$; \mathcal{E} 为等式约束指标集;

不等式约束: $c_i(\mathbf{x}) \geq 0, i \in \mathcal{J}$; \mathcal{J} 为不等式约束指标集

最优化理论数学基础

- 最优化问题:

$$\min f(\mathbf{x}),$$

$$\text{s. t. } \mathbf{x} \in \Omega$$

$$\Omega = \{\mathbf{x} \in R^n \mid c_i(\mathbf{x}) = 0, i \in \mathcal{E}; c_i(\mathbf{x}) \geq 0, i \in \mathcal{J}\}$$

无约束优化问题: $\mathcal{E} \cup \mathcal{J} = \emptyset$, 否则为约束优化问题

等式约束优化问题: $\mathcal{E} \neq \emptyset$ 且 $\mathcal{J} = \emptyset$

不等式约束优化问题: $\mathcal{J} \neq \emptyset$ 且 $\mathcal{E} = \emptyset$

线性规划问题: 目标函数和约束函数都是线性函数

二次规划问题: 目标函数二次, 而约束函数是线性

向量和矩阵范数

- 在 n 维实向量空间 R^n 中定义向量的范数
- 向量 $\mathbf{x} \in R^n$ 的范数 $\|\cdot\|$ 是一个非负数, 须满足以下条件:
 - ① $\|\mathbf{x}\| \geq 0, \|\mathbf{x}\| = 0 \Leftrightarrow \mathbf{x} = \mathbf{0}$
 - ② $\|\alpha\mathbf{x}\| = |\alpha|\|\mathbf{x}\|, \forall \alpha \in R$
 - ③ $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$
- p -范数: $\|\mathbf{x}\|_p = (\sum_{i=1}^n |x_i|^p)^{\frac{1}{p}}, p \geq 1$
 - 1-范数: $\|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i|$
 - 2-范数: $\|\mathbf{x}\|_2 = (\sum_{i=1}^n |x_i|^2)^{\frac{1}{2}}$
 - ∞ -范数: $\|\mathbf{x}\|_\infty = \max_{1 \leq i \leq n} |x_i|$

向量和矩阵范数

- 按照向量范数的定义引入矩阵 $A \in R^{m \times n}$ 的矩阵范数
- 矩阵范数 $\|\cdot\|$ 还需满足乘法性质:

$$\|AB\| \leq \|A\|\|B\|, \quad \forall A \in R^{m \times n}, B \in R^{n \times q}$$

- 称矩阵范数 $\|\cdot\|_\mu$ 与向量范数 $\|\cdot\|$ 是相容的, 当

$$\|Ax\| \leq \|A\|_\mu \|x\|, \quad \forall A \in R^{m \times n}, x \in R^n$$

- 若存在 $x \neq 0$ 使得

$$\|A\|_\mu = \max_{x \neq 0} \frac{\|Ax\|}{\|x\|} = \max_{\|x\|=1} \|Ax\|$$

称矩阵范数 $\|\cdot\|_\mu$ 是由向量范数 $\|\cdot\|$ 诱导出来的算子范数, 也称为从属于向量范数 $\|\cdot\|$ 的矩阵范数. 此时, 向量范数和算子范数通常采用相同的符号 $\|\cdot\|$.

向量和矩阵范数

- 显然, 从属于向量范数 $\|x\|_1$, $\|x\|_2$, $\|x\|_\infty$ 的矩阵范数为:

$$\|A\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^m |a_{ij}|$$

$$\|A\|_2 = \max\{\sqrt{\lambda} \mid \lambda \in \lambda(A^T A)\}$$

$$\|A\|_\infty = \max_{1 \leq i \leq m} \sum_{j=1}^n |a_{ij}|$$

分别称为列和范数, 谱范数, 行和范数.

注: $\lambda(A^T A)$ 表示 $A^T A$ 的特征值

向量和矩阵范数

- F-范数:

$$\|A\|_F = \left(\sum_{i=1}^m \sum_{j=1}^n a_{ij}^2 \right)^{1/2} = \sqrt{\text{tr}(A^T A)}$$

例如:

$$A = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \Rightarrow A^T = \begin{bmatrix} a_{11} & a_{21} \\ a_{12} & a_{22} \end{bmatrix}$$

$$A^T A = \begin{bmatrix} a_{11}^2 + a_{21}^2 & a_{11}a_{12} + a_{21}a_{22} \\ a_{11}a_{12} + a_{21}a_{22} & a_{12}^2 + a_{22}^2 \end{bmatrix}$$

$$\text{tr}(A^T A) = a_{11}^2 + a_{12}^2 + a_{21}^2 + a_{22}^2$$

谱范数和F-范数常用来讨论各种迭代算法的收敛性。

向量序列与矩阵序列的收敛性

- 若向量序列 $\{\mathbf{x}^{(k)}\}_{k=1}^{\infty} \subset \mathbf{R}^n$, 则:

$$\lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \mathbf{x} \Leftrightarrow \lim_{k \rightarrow \infty} x_i^{(k)} = x_i$$

- 类似, 若矩阵序列 $\{\mathbf{A}^{(k)}\}_{k=1}^{\infty} \subset \mathbf{R}^{m \times n}$, 则:

$$\lim_{k \rightarrow \infty} \mathbf{A}^{(k)} = \mathbf{A} \Leftrightarrow \lim_{k \rightarrow \infty} a_{ij}^{(k)} = a_{ij}$$

定理 1.1 (1) 设 $\|\cdot\|$ 和 $\|\cdot\|'$ 是定义在 \mathbf{R}^n 上的两个向量范数, 则存在两个正数 c_1 和 c_2 , 对所有 $\mathbf{x} \in \mathbf{R}^n$ 均成立

$$c_1 \|\mathbf{x}\| \leq \|\mathbf{x}\|' \leq c_2 \|\mathbf{x}\|.$$

(2) 设 $\|\cdot\|$ 和 $\|\cdot\|'$ 是定义在 $\mathbf{R}^{m \times n}$ 上的两个矩阵范数, 则存在两个正数 m_1 和 m_2 , 对所有 $\mathbf{A} \in \mathbf{R}^{m \times n}$ 均成立

$$m_1 \|\mathbf{A}\| \leq \|\mathbf{A}\|' \leq m_2 \|\mathbf{A}\|.$$

向量序列与矩阵序列的收敛性

定理 1.2 (1) 设 $\{\boldsymbol{x}^{(k)}\}$ 为 n 维向量序列, $\|\cdot\|$ 为定义在 \mathbf{R}^n 上的向量范数, 则

$$\lim_{k \rightarrow \infty} \boldsymbol{x}^{(k)} = \boldsymbol{x} \iff \lim_{k \rightarrow \infty} \|\boldsymbol{x}^{(k)} - \boldsymbol{x}\| = 0.$$

(2) 设 $\{\mathbf{A}^{(k)}\}$ 为 $m \times n$ 矩阵序列, $\|\cdot\|$ 为定义在 $\mathbf{R}^{m \times n}$ 上的矩阵范数, 则

$$\lim_{k \rightarrow \infty} \mathbf{A}^{(k)} = \mathbf{A} \iff \lim_{k \rightarrow \infty} \|\mathbf{A}^{(k)} - \mathbf{A}\| = 0.$$

多元函数分析

- 主要介绍三个概念
 - ✓ n 元函数的一阶导数
 - ✓ n 元函数的二阶导数
 - ✓ n 元函数的泰勒展开式

多元函数分析

- 定义1.1 $\mathbf{x} = (x_1, x_2, \dots, x_n)^T \in \mathbf{R}^n$. 称向量

$\nabla f(\mathbf{x}) = \left(\frac{\partial f(\mathbf{x})}{\partial x_1}, \frac{\partial f(\mathbf{x})}{\partial x_2}, \dots, \frac{\partial f(\mathbf{x})}{\partial x_n} \right)^T$ 为 $f(\mathbf{x})$ 在 \mathbf{x} 处的一阶导数或梯度. 称下列矩阵为 $f(\mathbf{x})$ 在 \mathbf{x} 处的二阶导数或Hesse阵

$$\nabla^2 f(\mathbf{x}) = \begin{pmatrix} \frac{\partial^2 f(\mathbf{x})}{\partial x_1^2} & \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_2 \partial x_1} & \frac{\partial^2 f(\mathbf{x})}{\partial x_2^2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_n \partial x_1} & \frac{\partial^2 f(\mathbf{x})}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_n^2} \end{pmatrix}$$

多元函数分析

- 若梯度 $\nabla f(\mathbf{x})$ 的每个分量函数在 \mathbf{x} 处都连续, 则称 f 在 \mathbf{x} 处**一阶连续可微**.
- 若Hesse阵 $\nabla^2 f(\mathbf{x})$ 的各个分量函数在 \mathbf{x} 处都连续, 则称 f 在 \mathbf{x} 处**二阶连续可微**.
- 若 f 在开集 \mathcal{D} 的每一点都一阶连续可微, 则称 f **在 \mathcal{D} 上一阶连续可微**.
- 若 f 在开集 \mathcal{D} 的每一点都二阶连续可微, 则称 f **在 \mathcal{D} 上二阶连续可微**.

多元函数分析

- 显然, 若 f 在 \mathbf{x} 处二阶连续可微, 则

$$\frac{\partial^2 f(\mathbf{x})}{\partial x_i \partial x_j} = \frac{\partial^2 f(\mathbf{x})}{\partial x_j \partial x_i}, \quad i, j = 1, 2, \dots, n$$

即 Hesse 阵 $\nabla^2 f(\mathbf{x})$ 是对称阵.

$$\nabla^2 f(\mathbf{x}) = \begin{pmatrix} \frac{\partial^2 f(\mathbf{x})}{\partial x_1^2} & \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_2 \partial x_1} & \frac{\partial^2 f(\mathbf{x})}{\partial x_2^2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_n \partial x_1} & \frac{\partial^2 f(\mathbf{x})}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_n^2} \end{pmatrix}$$

多元函数分析

- **泰勒展开:** 设函数 $f: \mathbb{R}^n \rightarrow \mathbb{R}$ 连续可微, 则

$$\begin{aligned} f(\mathbf{x} + \mathbf{h}) &= f(\mathbf{x}) + \int_0^1 \nabla f(\mathbf{x} + t\mathbf{h})^T \mathbf{h} dt \\ &= f(\mathbf{x}) + \nabla f(\mathbf{x} + \theta\mathbf{h})^T \mathbf{h} \quad (\theta \in (0,1)) \\ &= f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{h} + o(\|\mathbf{h}\|) \end{aligned}$$

$o(\|\mathbf{h}\|)$ 是泰勒展开的余项, 是一个一价无穷小

泰勒公式是一个用函数在某点的信息描述其附近取值的公式。如果函数足够平滑的话, 在已知函数在某一点的各阶导数值的情况之下, 泰勒公式可以用这些导数值做系数构建一个多项式来近似函数在这一点附近的邻域中的值。

$$f(x) = f(a) + f'(a)(x-a) + \frac{f''(a)}{2!}(x-a)^2 + \cdots + \frac{f^{(k)}(a)}{k!}(x-a)^k + h_k(x)(x-a)^k$$

多元函数分析

- **泰勒展开:** 设函数 $f: \mathbf{R}^n \rightarrow \mathbf{R}$ 连续可微, 则

$$f(\mathbf{x} + \mathbf{h}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{h} + o(\|\mathbf{h}\|)$$

若函数 f 是二次连续可微的, 则有

$$\begin{aligned} f(\mathbf{x} + \mathbf{h}) &= f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{h} + \int_0^1 (1-t) \mathbf{h}^T \nabla^2 f(\mathbf{x} + t\mathbf{h}) \mathbf{h} dt \\ &= f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{h} + \frac{1}{2} \mathbf{h}^T \nabla^2 f(\mathbf{x} + \theta\mathbf{h}) \mathbf{h} \quad (\theta \in (0,1)) \\ &= f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{h} + \frac{1}{2} \mathbf{h}^T \nabla^2 f(\mathbf{x}) \mathbf{h} + o(\|\mathbf{h}\|^2) \end{aligned}$$

拉格朗日余项

多元函数分析

- **泰勒展开:** 设函数 $f: \mathbf{R}^n \rightarrow \mathbf{R}$ 连续可微, 则

$$f(\mathbf{x} + \mathbf{h}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{h} + o(\|\mathbf{h}\|)$$

若函数 f 是二次连续可微的, 则有

$$f(\mathbf{x} + \mathbf{h}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{h} + \frac{1}{2} \mathbf{h}^T \nabla^2 f(\mathbf{x}) \mathbf{h} + o(\|\mathbf{h}\|^2)$$

以及

$$\begin{aligned} \nabla f(\mathbf{x} + \mathbf{h}) &= \nabla f(\mathbf{x}) + \int_0^1 \nabla^2 f(\mathbf{x} + t\mathbf{h})^T \mathbf{h} dt \\ &= \nabla f(\mathbf{x}) + \nabla^2 f(\mathbf{x} + \theta\mathbf{h})^T \mathbf{h} \quad (\theta \in (0,1)) \\ &= \nabla f(\mathbf{x}) + \nabla^2 f(\mathbf{x})^T \mathbf{h} + o(\|\mathbf{h}\|) \end{aligned}$$

多元函数分析

- 例1.1: 设二次函数

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T A \mathbf{x} - \mathbf{b}^T \mathbf{x}$$

式中: $\mathbf{b} \in \mathbb{R}^n$ 且 $A \in \mathbb{R}^{n \times n}$ 是**对称阵**. 求 $f(\mathbf{x})$ 在 \mathbf{x} 处的梯度和Hesse阵.

- 解: 显然 $\nabla f(\mathbf{x}) = A \mathbf{x} - \mathbf{b}$, $\nabla^2 f(\mathbf{x}) = A$

$$\text{例: } A = \begin{bmatrix} a_1 & a_2 \\ a_2 & a_3 \end{bmatrix}$$

对称阵条件实际中通常满足, 如: 距离矩阵等

多元函数分析

- 向量值函数

$$\mathbf{F}(\mathbf{x}) = (F_1(\mathbf{x}), F_2(\mathbf{x}), \dots, F_m(\mathbf{x}))^T: \mathbf{R}^n \rightarrow \mathbf{R}^m$$

- 若每个分量函数 F_i 都是(连续)可微的, 则称 \mathbf{F} 是(连续)可微的.

- 向量值函数 \mathbf{F} 在 \mathbf{x} 处的导数 $\mathbf{F}' \in \mathbf{R}^{m \times n}$ 称为它在 \mathbf{x} 处的**Jacobi矩阵**, 记为 $\mathbf{F}'(\mathbf{x})$ 或 $\mathbf{J}_F(\mathbf{x})$, 其转置称为

梯度矩阵 $\nabla \mathbf{F}(\mathbf{x}) = \mathbf{J}_F(\mathbf{x})^T$:

$$\mathbf{F}'(\mathbf{x}) := \mathbf{J}_F(\mathbf{x}) := \begin{pmatrix} \frac{\partial F_1(\mathbf{x})}{\partial x_1} & \frac{\partial F_1(\mathbf{x})}{\partial x_2} & \dots & \frac{\partial F_1(\mathbf{x})}{\partial x_n} \\ \frac{\partial F_2(\mathbf{x})}{\partial x_1} & \frac{\partial F_2(\mathbf{x})}{\partial x_2} & \dots & \frac{\partial F_2(\mathbf{x})}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial F_m(\mathbf{x})}{\partial x_1} & \frac{\partial F_m(\mathbf{x})}{\partial x_2} & \dots & \frac{\partial F_m(\mathbf{x})}{\partial x_n} \end{pmatrix}$$

多元函数分析

- 同理: 若向量值函数 $F: \mathbf{R}^n \rightarrow \mathbf{R}^m$ 是连续可微的, 则对于任意的 $\mathbf{x}, \mathbf{h} \in \mathbf{R}^n$, 有:

$$F(\mathbf{x} + \mathbf{h}) = F(\mathbf{x}) + F'(\mathbf{x})\mathbf{h} + o(\|\mathbf{h}\|)$$

定义 1.2 设向量值函数 $F: \mathbf{R}^n \rightarrow \mathbf{R}^m$, $\mathbf{x} \in \mathbf{R}^n$, 称 F 在 \mathbf{x} 处是 Lipschitz 连续的, 是指存在常数 $L > 0$, 使得对任意的 $\mathbf{y} \in \mathbf{R}^n$, 满足

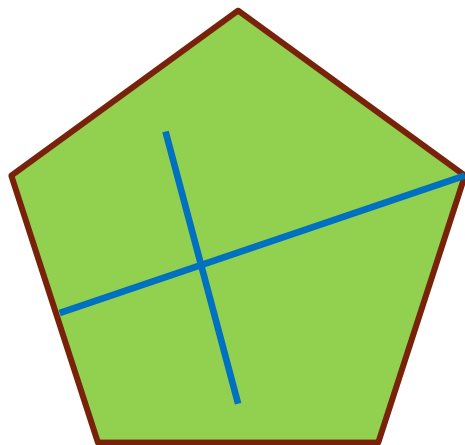
$$\|F(\mathbf{x}) - F(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|, \quad (1.8)$$

式中: L 为 Lipschitz 常数.

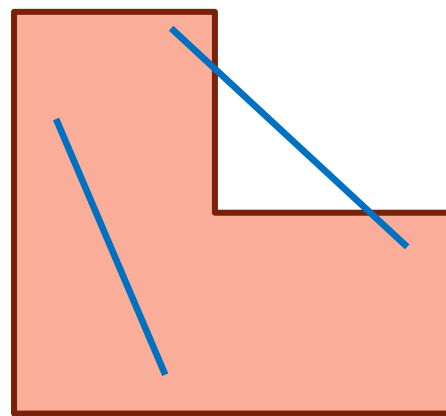
若式 (1.8) 对任意的 $\mathbf{x}, \mathbf{y} \in \mathbf{R}^n$ 都成立, 则称 F 在 \mathbf{R}^n 内是 Lipschitz 连续的.

凸集与凸函数

- 介绍凸集、锥和凸函数的有关概念
- **定义1.3** 设集合 $\mathcal{D} \subset \mathbb{R}^n$. 称集合 \mathcal{D} 为**凸集**, 是指对任意的 $x, y \in \mathcal{D}$ 及任意的实数 $\lambda \in [0, 1]$, 都有 $\lambda x + (1 - \lambda)y \in \mathcal{D}$.



凸集



非凸

凸集与凸函数

- **性质1.1** 设 $\mathcal{D}, \mathcal{D}_1, \mathcal{D}_2$ 是凸集, α 是一实数, 那么
 - (1) $\alpha\mathcal{D} := \{y | y = \alpha x, x \in \mathcal{D}\}$ 是凸集
 - (2) **交集** $\mathcal{D}_1 \cap \mathcal{D}_2$ 是凸集
 - (3) **和集** $\mathcal{D}_1 + \mathcal{D}_2 := \{z | z = x + y, x \in \mathcal{D}_1, y \in \mathcal{D}_2\}$ 是凸集
- **例1.3** n 维欧几里得空间中的 m 个点的凸组合是一个凸集, 即集合

$$\left\{ x = \sum_{i=1}^m \alpha_i x_i \mid x_i \in R^n, \alpha_i \geq 0, \sum_{i=1}^m \alpha_i = 1 \right\}$$

是凸集

凸集与凸函数

- 定义1.4 集合 $D \subset \mathbb{R}^n$ 的**凸包(convex hull)**是指所有包含 D 的凸集的交集, 记为

$$\text{conv}(D) := \bigcap_{G \supseteq D} G$$

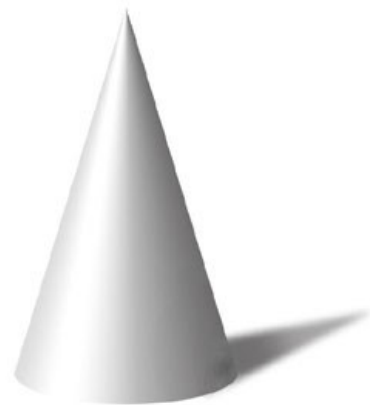
式中: G 为凸集



凸集与凸函数

- **定义1.5** 设非空集合 $G \subset R^n$. 若对任意的 $x \in G$ 和任意的实数 $\lambda > 0$, 有 $\lambda x \in G$, 则称 G 为一个**锥 (cone)**. 若 G 同时也是凸集, 则称 G 为一个**凸锥 (convex cone)**. 此外, 对于锥 G , 若 $0 \in G$, 则称 G 为一个**尖锥 (pointed cone)**. 相应地, 包含 0 的凸锥成为**尖凸锥**.

- **例1.6** 多面体 $\{x \in R^n | Ax \geq 0\}$ 是一个尖凸锥, 通常称为多面锥



凸集与凸函数

- 定义凸集上的凸函数

- **定义1.6** 设函数 $f: \mathcal{D} \subset \mathbb{R}^n \rightarrow \mathbb{R}$, 其中 \mathcal{D} 为凸集.

(1) 称 f 是 \mathcal{D} 上的 **凸函数**, 是指对任意的 $\mathbf{x}, \mathbf{y} \in \mathcal{D}$ 及任意的实数 $\lambda \in [0, 1]$, 都有

$$f(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda) f(\mathbf{y})$$

(2) 称 f 是 \mathcal{D} 上的 **严格凸函数**, 是指对任意的 $\mathbf{x}, \mathbf{y} \in \mathcal{D}$ 及任意的实数 $\lambda \in (0, 1)$, 都有

$$f(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) < \lambda f(\mathbf{x}) + (1 - \lambda) f(\mathbf{y})$$

(3) 称 f 是 \mathcal{D} 上的 **一致凸函数**, 是指存在常数 $\gamma > 0$, 使对任意的 $\mathbf{x}, \mathbf{y} \in \mathcal{D}$ 及任意的实数 $\lambda \in [0, 1]$, 都有

$$f(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) + \frac{1}{2} \lambda (1 - \lambda) \gamma \|\mathbf{x} - \mathbf{y}\|^2 \leq \lambda f(\mathbf{x}) + (1 - \lambda) f(\mathbf{y})$$

凸集与凸函数

- 定义凸集上的凸函数

- **性质1.2** 设 f, f_1, f_2 都是凸集 \mathcal{D} 上的凸函数, $c_1, c_2 \in \mathbf{R}_+$, $\alpha \in \mathbf{R}$, 则有

- (1) $c_1 f_1(\mathbf{x}) + c_2 f_2(\mathbf{x})$ 也是 \mathcal{D} 上的凸函数;

- (2) 水平集 $\mathcal{L}(f, \alpha) = \{\mathbf{x} | \mathbf{x} \in \mathcal{D}, f(\mathbf{x}) \leq \alpha\}$ 是凸集.

- 凸集和凸函数在传统优化理论中起着举足轻重的作用, 但是用凸函数的定义来判断一个函数是否具有凸性相当困难.
- 如果函数是一阶或二阶连续可微的, 那么利用函数的梯度或Hesse矩阵来判别或验证函数的凸性要相对容易.

函数的凸性

● 定理 1.4 设 f 在凸集 $\mathcal{D} \subset \mathbb{R}^n$ 上一阶连续可微, 则

(1) f 在 \mathcal{D} 上为凸函数的充要条件是

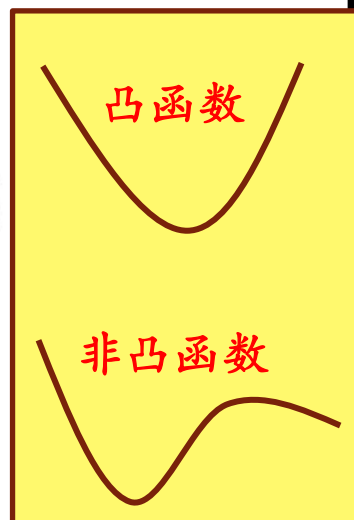
$$f(\mathbf{x}) \geq f(\mathbf{y}) + \nabla f(\mathbf{y})^T(\mathbf{x} - \mathbf{y}), \quad \forall \mathbf{x}, \mathbf{y} \in \mathcal{D};$$

(2) f 在 \mathcal{D} 上为严格凸函数的充要条件是, 当 $\mathbf{x} \neq \mathbf{y}$ 时, 成立

$$f(\mathbf{x}) > f(\mathbf{y}) + \nabla f(\mathbf{y})^T(\mathbf{x} - \mathbf{y}), \quad \forall \mathbf{x}, \mathbf{y} \in \mathcal{D};$$

(3) f 在 \mathcal{D} 上为一致凸函数的充要条件是, 存在常数 $c > 0$, 对任意的 $\mathbf{x}, \mathbf{y} \in \mathcal{D}$, 成立

$$f(\mathbf{x}) \geq f(\mathbf{y}) + \nabla f(\mathbf{y})^T(\mathbf{x} - \mathbf{y}) + c\|\mathbf{x} - \mathbf{y}\|^2.$$



函数的凸性

- 在一元函数中, 若 $f(x)$ 在区间 (a, b) 上二阶可微且 $f''(x) \geq 0 (> 0)$, 则 $f(x)$ 在 (a, b) 内凸(严格凸).
- 对于二阶连续可微的多元函数 $f: D \subset \mathbb{R}^n \rightarrow \mathbb{R}$, 也可以由其二阶导数(Hesse矩阵)给出凸性的一个近乎完整的表述.

定义 1.7 设 n 元实函数 f 在凸集 D 上是二阶连续可微的. 若对一切 $h \in \mathbb{R}^n$, 有 $h^T \nabla^2 f(x) h \geq 0$, 则称 $\nabla^2 f$ 在点 x 处是半正定的. 若对一切 $0 \neq h \in \mathbb{R}^n$, 有 $h^T \nabla^2 f(x) h > 0$, 则称 $\nabla^2 f$ 在点 x 处是正定的. 进一步, 若存在常数 $c > 0$, 使得对任意的 $h \in \mathbb{R}^n$, $x \in D$, 有 $h^T \nabla^2 f(x) h \geq c \|h\|^2$, 则称 $\nabla^2 f$ 在 D 上是一致正定的.

函数的凸性

- 在一元函数中, 若 $f(x)$ 在区间 (a, b) 上二阶可微且 $f''(x) \geq 0 (> 0)$, 则 $f(x)$ 在 (a, b) 内凸(严格凸).
- 推广到多元函数
- **定理1.5** 设 n 元实函数 f 在凸集 $D \subset R^n$ 上二阶连续可微, 则:
 - (1) f 在 D 上为凸函数的充要条件是 $\nabla^2 f(x)$ 对一切 $x \in D$ 为半正定;
 - (2) f 在 D 上为严格凸函数的充分条件是 $\nabla^2 f(x)$ 对一切 $x \in D$ 为正定;
 - (3) f 在 D 上为一致凸函数的充要条件是 $\nabla^2 f(x)$ 对一切 $x \in D$ 为一致正定.

无约束问题的最优性条件

- 讨论无约束优化问题

$$\min_{x \in \mathbb{R}^n} f(x)$$

的**最优性条件**

- ✓ 一阶必要条件
- ✓ 二阶必要条件
- ✓ 二阶充分条件
- ✓ 凸函数全局极小点的充要条件

无约束问题的最优性条件

- 定义 1.8 若对于任意的 $x \in \mathbf{R}^n$, 都有

$$f(x^*) \leq f(x),$$

则称 x^* 为 f 的一个全局极小点. 若上述不等式严格成立且 $x \neq x^*$, 则称 x^* 为 f 的一个严格全局极小点.

- 定义 1.9 若对于任意的 $x \in N(x^*, \delta) = \{x \in \mathbf{R}^n \mid \|x - x^*\| < \delta\}$, 都有

$$f(x^*) \leq f(x),$$

则称 x^* 为 f 的一个局部极小点, 其中 $\delta > 0$ 为某个常数. 若上述不等式严格成立且 $x \neq x^*$, 则称 x^* 为 f 的一个严格局部极小点.

无约束问题的最优性条件

- 一般讨论的求极小点的方法都是指局部极小点.
- 为了表述的方便, 引入下列记号:

$$\mathbf{g}(\mathbf{x}) = \nabla f(\mathbf{x}), \mathbf{g}_k = \nabla f(\mathbf{x}_k)$$

$$\mathbf{G}(\mathbf{x}) = \nabla^2 f(\mathbf{x}), \mathbf{G}_k = \nabla^2 f(\mathbf{x}_k)$$

无约束问题的最优性条件

- **定理1.6** (一阶必要条件) 设 $f(\mathbf{x})$ 在开集 \mathcal{D} 上一阶连续可微. 若 $\mathbf{x}^* \in \mathcal{D}$ 是问题 $\min_{\mathbf{x} \in \mathcal{D}} f(\mathbf{x})$ 的一个局部极小点, 则必有 $\mathbf{g}(\mathbf{x}^*) = \mathbf{0}$.

证明 取 $\mathbf{x} = \mathbf{x}^* - \alpha \mathbf{g}(\mathbf{x}^*) \in \mathcal{D}$, 其中 $\alpha > 0$ 为某个常数, 则有

$$\begin{aligned} f(\mathbf{x}) &= f(\mathbf{x}^*) + \mathbf{g}(\mathbf{x}^*)^\top (\mathbf{x} - \mathbf{x}^*) + o(\|\mathbf{x} - \mathbf{x}^*\|) \\ &= f(\mathbf{x}^*) - \alpha \mathbf{g}(\mathbf{x}^*)^\top \mathbf{g}(\mathbf{x}^*) + o(\alpha) \\ &= f(\mathbf{x}^*) - \alpha \|\mathbf{g}(\mathbf{x}^*)\|^2 + o(\alpha). \end{aligned}$$

注意到 $f(\mathbf{x}) \geq f(\mathbf{x}^*)$ 及 $\alpha > 0$, 于是有

$$0 \leq \|\mathbf{g}(\mathbf{x}^*)\|^2 \leq \frac{o(\alpha)}{\alpha}.$$

上式两边令 $\alpha \rightarrow 0$, 得 $\|\mathbf{g}(\mathbf{x}^*)\| = 0$, 即 $\mathbf{g}(\mathbf{x}^*) = \mathbf{0}$. 证毕.

无约束问题的最优性条件

- **定理1.7** (二阶必要条件) 设 $f(\mathbf{x})$ 在开集 \mathcal{D} 上二阶连续可微. 若 $\mathbf{x}^* \in \mathcal{D}$ 是问题 $\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$ 的一个局部极小点, 则必有 $\mathbf{g}(\mathbf{x}^*) = 0$ 且 $\mathbf{G}(\mathbf{x}^*)$ 是 **半正定矩阵**.

证明 设 \mathbf{x}^* 是一局部极小点, 那么由定理 1.6 可知 $\mathbf{g}(\mathbf{x}^*) = 0$. 下面只需证明 $\mathbf{G}(\mathbf{x}^*)$ 的半正定性. 任取 $\mathbf{x} = \mathbf{x}^* + \alpha \mathbf{d} \in \mathcal{D}$, 其中 $\alpha > 0$

且 $\mathbf{d} \in \mathbb{R}^n$. 由泰勒展开式, 得 $f(\mathbf{x} + \mathbf{h}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{h} + \frac{1}{2} \mathbf{h}^T \nabla^2 f(\mathbf{x}) \mathbf{h} + o(\|\mathbf{h}\|^2)$

$$0 \leq f(\mathbf{x}) - f(\mathbf{x}^*) = \frac{1}{2} \alpha^2 \mathbf{d}^T \mathbf{G}(\mathbf{x}^*) \mathbf{d} + o(\alpha^2),$$

即

$$\mathbf{d}^T \mathbf{G}(\mathbf{x}^*) \mathbf{d} + \frac{o(2\alpha^2)}{\alpha^2} \geq 0.$$

对上式令 $\alpha \rightarrow 0$, 即得 $\mathbf{d}^T \mathbf{G}(\mathbf{x}^*) \mathbf{d} \geq 0$, 从而定理成立. 证毕.

无约束问题的最优性条件

- **定理1.8** (二阶充分条件) 设 $f(x)$ 在开集 \mathcal{D} 上二阶连续可微. 若 $x^* \in \mathcal{D}$ 满足条件 $g(x^*) = 0$ 且 $G(x^*)$ 是正定矩阵, 则 x^* 是问题 $\min_{x \in \mathbb{R}^n} f(x)$ 的一个局部极小点.
- 目标函数的稳定点(驻点, 一阶导数为0)不一定是极小点.
- 但是, 目标函数是凸函数的无约束优化问题, 其稳定点、局部极小点和全局极小点等价.

无约束问题的最优性条件

- 定理1.9 设 $f(x)$ 在 R^n 上是凸函数并且一阶连续可微. 则 $x^* \in R^n$ 是问题 $\min_{x \in R^n} f(x)$ 的全局极小点的充要条件是 $g(x^*) = 0$.

证明 必要性是显然的, 所以只需证明充分性. 设 $g(x^*) = 0$. 由凸函数的判别定理 1.4 (1), 得

$$f(x) \geq f(x^*) + g(x^*)^T(x - x^*) = f(x^*), \quad \forall x \in R^n,$$

f 在 \mathcal{D} 上为凸函数的充要条件是

$$f(x) \geq f(y) + \nabla f(y)^T(x - y), \quad \forall x, y \in \mathcal{D};$$

无约束优化问题的数值优化方法

- 在数值优化中,一般采用迭代法求解无约束优化问题的极小点.
- 基本思路如下:
 - ✓ 给定一初始点 x_0 ;
 - ✓ 按照某迭代规则产生一迭代序列 $\{x_k\}$,使该序列是**有限的**,则**最后一个点**就是极小点;
 - ✓ 若序列 $\{x_k\}$ 是**无限的**,它有**极限点**且该点即为极小点.
- 迭代细节:
 - x_k 为第 k 次迭代点, d_k 为第 k 次**搜索方向**, α_k 为第 k 次**步长因子**,则可得第 $k+1$ 次迭代点 $x_{k+1} = x_k + \alpha_k d_k$.

无约束优化问题的数值优化方法

● 无约束优化问题的一般算法框架

算法 1.1 (无约束优化问题的一般算法框架)

步骤 0, 给定初始化参数及初始迭代点 \mathbf{x}_0 . 令 $k := 0$.

步骤 1, 若 \mathbf{x}_k 满足某种终止准则, 停止迭代, 以 \mathbf{x}_k 作为近似极小点.

步骤 2, 通过求解 \mathbf{x}_k 处的某个子问题确定下降方向 \mathbf{d}_k .

步骤 3, 通过某种搜索方式确定步长因子 α_k , 使得 $f(\mathbf{x}_k + \alpha_k \mathbf{d}_k) < f(\mathbf{x}_k)$

步骤 4, 令 $\mathbf{x}_{k+1} := \mathbf{x}_k + \alpha_k \mathbf{d}_k$, $k := k + 1$, 转步骤 1.

- 称上述算法中的 $\mathbf{s}_k = \alpha_k \mathbf{d}_k$ 为第 k 次迭代的位移.
- 设计不同的位移(不同搜索方向和步长因子)会产生不同的迭代算法. 为了保证收敛性, 一般要求搜索方向为所谓的下降方向.

无约束优化问题的数值优化方法

- **定义 1.10** 若存在 $\bar{\alpha} > 0$, 使得对任意的 $\alpha \in (0, \bar{\alpha})$ 和 $\mathbf{d}_k \neq \mathbf{0}$, 有 $f(\mathbf{x}_k + \alpha \mathbf{d}_k) < f(\mathbf{x}_k)$, 则称 \mathbf{d}_k 为 $f(\mathbf{x})$ 在 \mathbf{x}_k 处的一个 **下降方向**.
- **引理 1.1** 设函数 $f: \mathcal{D} \subset \mathbf{R}^n \rightarrow \mathbf{R}$ 在开集 \mathcal{D} 上一阶连续可微. 则 \mathbf{d}_k 为 $f(\mathbf{x})$ 在 \mathbf{x}_k 处的一个 **下降方向** 的充要条件是.

$$\nabla f(\mathbf{x}_k)^T \mathbf{d}_k < 0$$

$$f(\mathbf{x}_k + \alpha \mathbf{d}_k) = f(\mathbf{x}_k) + \alpha \nabla f(\mathbf{x}_k)^T \mathbf{d}_k + o(\alpha).$$

无约束优化问题的数值优化方法

- **定义 1.11** 若某算法只有当初始点 x_0 充分接近极小点 x^* 时, 由算法产生的点列 $\{x_k\}$ 才收敛于 x^* , 则称该算法具有**局部收敛性**. 若对于任意的初始点 x_0 , 由算法产生的点列 $\{x_k\}$ 都收敛于 x^* , 则称该算法具有**全局收敛性**.
- 算法的局部收敛速度是衡量一个算法好坏的重要指标, 通常有两种衡量收敛速度的尺度: Q-收敛和R-收敛.

无约束优化问题的数值优化方法

● Q-收敛

定义 1.12 设算法产生的点列 $\{\mathbf{x}_k\}$ 收敛于极小点 \mathbf{x}^* , 且

$$\limsup_{k \rightarrow \infty} \frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}_k - \mathbf{x}^*\|^p} \leq c.$$

(1) 若 $p = 1$ 且 $0 < c < 1$, 则称该算法具有 Q-线性收敛速度 (或 Q-线性收敛的).

(2) 若 $p = 1$ 且 $c = 0$, 则称该算法具有 Q-超线性收敛速度 (或 Q-超线性收敛的).

(3) 若 $p = 2$ 且 $0 < c < \infty$, 则称该算法具有 Q-平方收敛速度 (或 Q-平方收敛的).

(4) 一般地, 若 $p > 2$ 且 $0 < c < \infty$, 则称该算法具有 Q- p 阶收敛速度 (或 Q- p 阶收敛的).

无约束优化问题的数值优化方法

● R-收敛

定义 1.13 设算法产生的点列 $\{\mathbf{x}_k\}$ 收敛于极小点 \mathbf{x}^* .

(1) 若存在常数 $c > 0$ 和 $q \in (0, 1)$ 使得

$$\|\mathbf{x}_k - \mathbf{x}^*\| \leq cq^k,$$

则称序列 $\{\mathbf{x}_k\}$ R-线性收敛到 \mathbf{x}^* .

(2) 若存在常数 $c > 0$ 和收敛于零的正数列 $\{q_k\}$ 使得

$$\|\mathbf{x}_k - \mathbf{x}^*\| \leq c \prod_{i=0}^k q_i,$$

则称序列 $\{\mathbf{x}_k\}$ R-超线性收敛到 \mathbf{x}^* .

无约束优化问题的数值优化方法

● 迭代算法常用的终止条件:

(1) 位移的绝对误差或相对误差充分小, 即

$$\|\mathbf{x}_{k+1} - \mathbf{x}_k\| < \varepsilon \quad \text{或} \quad \frac{\|\mathbf{x}_{k+1} - \mathbf{x}_k\|}{\|\mathbf{x}_k\|} < \varepsilon,$$

式中: ε 为充分小的正数;

(2) 目标函数的绝对误差或相对误差充分小, 即

$$|f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k)| < \varepsilon \quad \text{或} \quad \frac{|f(\mathbf{x}_{k+1}) - f(\mathbf{x}_k)|}{|f(\mathbf{x}_k)|} < \varepsilon,$$

式中: ε 为充分小的正数;

(3) 目标函数的梯度的范数充分小, 即

$$\|\nabla f(\mathbf{x}_k)\| < \varepsilon,$$

式中: ε 为充分小的正数.

线搜索方法

- 线搜索方法是求解无约束优化问题 $\min_{x \in \mathbb{R}^n} f(x)$ 的一个最基本的方法. 此处介绍 **一维线搜索** 方法.

算法 1.1 (无约束优化问题的一般算法框架)

步骤 0, 给定初始化参数及初始迭代点 x_0 . 令 $k := 0$.

步骤 1, 若 x_k 满足某种终止准则, 停止迭代, 以 x_k 作为近似极小点.

步骤 2, 通过求解 x_k 处的某个子问题确定下降方向 d_k .

步骤 3, 通过某种搜索方式确定步长因子 α_k , 使得 $f(x_k + \alpha_k d_k) < f(x_k)$

步骤 4, 令 $x_{k+1} := x_k + \alpha_k d_k$, $k := k + 1$, 转步骤 1.

线搜索方法

- 线搜索方法是求解无约束优化问题 $\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$ 的一个最基本的方法. 此处介绍 **一维线搜索** 方法.

步骤 3, 通过某种搜索方式确定步长因子 α_k , 使得 $f(\mathbf{x}_k + \alpha_k \mathbf{d}_k) < f(\mathbf{x}_k)$

- 这实际上是将原优化问题转变为(n 个变量)目标函数 $f(\mathbf{x})$ 在一个规定方向上移动所形成的单变量优化问题, 也即所谓的“**线搜索**”或“**一维搜索**”. 令

$$\phi(\alpha) = f(\mathbf{x}_k + \alpha \mathbf{d}_k)$$

如此, 搜索步骤3等价于求步长因子 α_k 使得 $\phi(\alpha_k) < \phi(0)$.

线搜索方法

$$\phi(\alpha) = f(\mathbf{x}_k + \alpha \mathbf{d}_k)$$

求步长因子 α_k 使得 $\phi(\alpha_k) < \phi(0)$.

- **精确线搜索**: 求 α_k 使目标函数 f 沿方向 \mathbf{d}_k 达到极小, 即

$$\phi(\alpha_k) = \min_{\alpha > 0} \phi(\alpha)$$

若 $f(x)$ 连续可微, 那么步长因子 α_k 具有如下性质:

$$\nabla f(\mathbf{x}_k + \alpha_k \mathbf{d}_k)^T \mathbf{d}_k = 0 \quad (\text{也即 } \mathbf{g}_{k+1}^T \mathbf{d}_k = 0)$$

- **非精确线搜索**: 求 α_k 使目标函数 f 得到可接受的下降量, 即

$$\Delta f_k = f(\mathbf{x}_k) - f(\mathbf{x}_k + \alpha_k \mathbf{d}_k) > 0$$

是可接受的.

精确线搜索方法

- **精确线搜索基本思想**: 首先**确定**包含问题最优解的搜索**区间**, 然后采用某种插值或分割技术**缩小**这个**区间**, 进行搜索求解.
- **搜索区间定义**:

定义 2.1 设 ϕ 是定义在实数集上一元实函数, $\alpha^* \in [0, +\infty)$, 并且

$$\phi(\alpha^*) = \min_{\alpha \geq 0} \phi(\alpha). \quad (2.5)$$

若存在区间 $[a, b] \subset [0, +\infty)$, 使 $\alpha^* \in (a, b)$, 则称 $[a, b]$ 是极小化问题 (2.5) 的**搜索区间**. 进一步, 若 α^* 使得 $\phi(\alpha)$ 在 $[a, \alpha^*]$ 上严格递减, 在 $[\alpha^*, b]$ 上严格递增, 则称 $[a, b]$ 是 $\phi(\alpha)$ 的**单峰区间**, $\phi(\alpha)$ 是 $[a, b]$ 上的**单峰函数**.

精确线搜索方法

- **精确线搜索**: 首先**确定**包含问题最优解的搜索**区间**, 然后采用某种插值或分割技术**缩小**这个**区间**, 进行搜索求解.
- ① 如何确定搜索区间并保证具有近似单峰性质?
- **进退法基本思想**: 从一点出发, 按一定步长, 试图确定函数值呈现“**高-低-高**”的三点, 从而得到一个近似的单峰区间.

精确线搜索方法

① 如何确定搜索区间并保证具有近似单峰性质?

● 进退法:

算法 2.1 (确定近似单峰区间的进退法)

步骤 0, 选取 $\alpha_0 \geq 0$, $h_0 > 0$. 计算 $\phi_0 := \phi(\alpha_0)$. 令 $k := 0$.

步骤 1, 令 $\alpha_{k+1} := \alpha_k + h_k$, 计算 $\phi_{k+1} := \phi(\alpha_{k+1})$. 若 $\phi_{k+1} < \phi_k$, 转步骤 2; 否则, 转步骤 3.

步骤 2, 加大步长. 令 $h_{k+1} := 2h_k$, $\alpha := \alpha_k$, $\alpha_k := \alpha_{k+1}$, $\phi_k := \phi_{k+1}$, $k := k + 1$, 转步骤 1.

步骤 3, 反向搜索或输出. 若 $k = 0$, 令 $h_1 := -h_0$, $\alpha := \alpha_1$, $\alpha_1 := \alpha_0$, $\phi_1 := \phi_0$, $k := 1$, 转步骤 1; 否则, 停止迭代, 令

$$a = \min\{\alpha, \alpha_{k+1}\}, \quad b = \max\{\alpha, \alpha_{k+1}\}.$$

输出 $[a, b]$.

α 实际上相当于是 α_{k-1}

精确线搜索方法

- **精确线搜索**: 首先**确定**包含问题最优解的搜索**区间**, 然后采用某种插值或分割技术**缩小**这个**区间**, 进行搜索求解.
- ② 如何采用某种插值或分割技术缩小这个区间?
- **方法一般分为两类**: 1) 使用导数的搜索, 如插值法、牛顿法及**抛物线法**等; 2) 不使用导数的搜索, 如**黄金分割法**、分数法等.
- **黄金分割法基本思想**: 通过试探点函数值的比较, 使包含极值点的搜索区间不断缩小. 该方法仅需要计算函数值.
- **抛物线法基本思想**: 通过在搜索区间中不断地使用二次多项式去近似目标函数, 并逐步用插值多项式的极小点去逼近线搜索问题的极小点. 该方法也不需要计算导数值.

精确线搜索方法——黄金分割法

算法 2.2 (黄金分割法)

步骤 0, 确定初始搜索区间 $[a_0, b_0]$ 和容许误差 $0 \leq \varepsilon \ll 1$. 令

$t = (\sqrt{5} - 1)/2$, 计算初始试探点

$$p_0 = a_0 + (1 - t)(b_0 - a_0), \quad q_0 = a_0 + t(b_0 - a_0)$$

及相应的函数值 $\phi(p_0), \phi(q_0)$. 令 $i := 0$.

步骤 1, 若 $\phi(p_i) \leq \phi(q_i)$, 转步骤 2; 否则, 转步骤 3.

步骤 2, 计算左试探点. 若 $|q_i - a_i| \leq \varepsilon$, 停算, 输出 p_i 否则, 令

$$a_{i+1} := a_i, \quad b_{i+1} := q_i, \quad \phi(q_{i+1}) := \phi(p_i),$$

$$q_{i+1} := p_i, \quad p_{i+1} := a_{i+1} + (1 - t)(b_{i+1} - a_{i+1}).$$

计算 $\phi(p_{i+1})$, $i := i + 1$, 转步骤 1.

步骤 3, 计算右试探点. 若 $|b_i - p_i| \leq \varepsilon$, 停算, 输出 q_i 否则, 令

$$a_{i+1} := p_i, \quad b_{i+1} := b_i, \quad \phi(p_{i+1}) := \phi(q_i),$$

$$p_{i+1} := q_i, \quad q_{i+1} := a_{i+1} + t(b_{i+1} - a_{i+1}).$$

计算 $\phi(q_{i+1})$, $i := i + 1$, 转步骤 1.

Why?
t为何取
该值?

保持高低高
结构

精确线搜索方法——黄金分割法

- 第 i 次迭代时, 搜索区间为 $[a_i, b_i]$. 取两个试探点为 $p_i, q_i \in [a_i, b_i]$ 且 $p_i < q_i$, 并满足以下两个条件:
 - ✓ $[a_i, q_i]$ 与 $[p_i, b_i]$ 的**长度相同**, 即 $b_i - p_i = q_i - a_i$
 - ✓ 区间长度的**缩短率相同**, 即 $b_{i+1} - a_{i+1} = t(b_i - a_i)$

黄金分割法每次迭代搜索区间搜索率是0.618, 只是线性收敛, 计算效率不高, 但是每次迭代只需计算一次函数值.

$$p_i = a_i + (1 - t)(b_i - a_i), \quad q_i = a_i + t(b_i - a_i).$$

考虑情形 $\phi(p_i) \leq \phi(q_i)$, 新的搜索区间为

$$[a_{i+1}, b_{i+1}] = [a_i, q_i].$$

$$\begin{aligned} q_{i+1} &= a_{i+1} + t(b_{i+1} - a_{i+1}) = a_i + t(q_i - a_i) \\ &= a_i + t^2(b_i - a_i). \end{aligned}$$

若令

$$t^2 = 1 - t, \quad t > 0,$$

则 $q_{i+1} = a_i + (1 - t)(b_i - a_i) = p_i.$

新的试探点 q_{i+1} 就不需要重新计算.

得区间长度缩短率为 $t = \frac{\sqrt{5} - 1}{2} \approx 0.618.$

精确线搜索方法——抛物线法

- **抛物线法基本思想**: 通过在搜索区间中不断地使用二次多项式去近似目标函数, 并逐步用插值多项式的极小点去逼近线搜索问题的极小点.

算法 2.3 (抛物线法)

步骤 0, 由算法 2.1 确定三点 $s_0 < s_1 < s_2$, 对应的函数值

ϕ_0, ϕ_1, ϕ_2 满足 $\phi_1 < \phi_0, \phi_1 < \phi_2$.

设定容许误差 $0 \leq \varepsilon \ll 1$.

$$\bar{s} = s_0 + \frac{(3\phi_0 - 4\phi_1 + \phi_2)h}{2(\phi_0 - 2\phi_1 + \phi_2)}$$

步骤 1, 若 $|s_2 - s_0| < \varepsilon$, 停算, 输出 $s^* \approx s_1$.

步骤 2, 计算插值点. 根据式 (2.9) 计算 \bar{s} 和 $\bar{\phi} := \phi(\bar{s})$. 若 $\phi_1 \leq \bar{\phi}$, 转步骤 4; 否则, 转步骤 3.

步骤 3, 若 $s_1 > \bar{s}$, 则 $s_2 := s_1, s_1 := \bar{s}, \phi_2 := \phi_1, \phi_1 := \bar{\phi}$, 转步骤 1; 否则, $s_0 := s_1, s_1 := \bar{s}, \phi_0 := \phi_1, \phi_1 := \bar{\phi}$, 转步骤 1.

步骤 4, 若 $s_1 < \bar{s}$, 则 $s_2 := \bar{s}, \phi_2 := \bar{\phi}$, 转步骤 1; 否则, $s_0 := \bar{s}, \phi_0 := \bar{\phi}$, 转步骤 1.

精确线搜索方法——抛物线法

- 抛物线法原理解释: 1)如何得到 \bar{s} ? 2)搜索的最优性保证?

已知三点 $s_0, s_1 = s_0 + h, s_2 = s_0 + 2h$ ($h > 0$) 处的函数值 ϕ_0, ϕ_1, ϕ_2 , 且满足 $\phi_1 < \phi_0, \phi_1 < \phi_2$. 保证了函数 ϕ 在区间 $[s_0, s_2]$ 上是单峰函数. 则 **三点间二次拉格朗日插值**

多项式为

$$q(s) = \frac{(s - s_1)(s - s_2)}{2h^2}\phi_0 - \frac{(s - s_0)(s - s_2)}{h^2}\phi_1 + \frac{(s - s_0)(s - s_1)}{2h^2}\phi_2.$$

$q(s)$ 的一阶导数为

$$q'(s) = \frac{2s - s_1 - s_2}{2h^2}\phi_0 - \frac{2s - s_0 - s_2}{h^2}\phi_1 + \frac{2s - s_0 - s_1}{2h^2}\phi_2.$$

令 $q'(s) = 0$, 解得

$$\begin{aligned}\bar{s} &= \frac{(s_1 + s_2)\phi_0 - 2(s_0 + s_2)\phi_1 + (s_0 + s_1)\phi_2}{2(\phi_0 - 2\phi_1 + \phi_2)} \\ &= \frac{(2s_0 + 3h)\phi_0 - 2(2s_0 + 2h)\phi_1 + (2s_0 + h)\phi_2}{2(\phi_0 - 2\phi_1 + \phi_2)} \\ &= s_0 + \frac{(3\phi_0 - 4\phi_1 + \phi_2)h}{2(\phi_0 - 2\phi_1 + \phi_2)} := s_0 + \bar{h},\end{aligned}$$

因 $q(s)$ 的二阶导数为

$$q''(s) = \frac{\phi_0}{h^2} - \frac{2\phi_1}{h^2} + \frac{\phi_2}{h^2} = \frac{\phi_0 - 2\phi_1 + \phi_2}{h^2} > 0,$$

问: 为何不直接解 $\phi'(s) = 0$?

拟合三点间的二次曲线

求拟合曲线的极小值

凸二次连续可微, 保证最优性

(2.9)

故 $q(s)$ 为凸二次函数.

s_{\min} 是 $q(s)$ 的全局极小点.

非精确线搜索方法

- **精确线搜索缺点**: 往往需要计算很多的函数值和梯度值, 特别当迭代点远离最优点时, 不是十分有效合理.
- **非精确线搜索**: 既能保证目标函数具有可接受的下降量, 又能使最终形成的迭代序列收敛.
 - ✓ Wolfe 准则
 - ✓ Armijo 准则

非精确线搜索方法——Wolfe准则

- Wolfe准则是指给定 $\rho \in (0, 0.5)$, $\sigma \in (\rho, 1)$, 求使得下面两个不等式同时成立:

$$f(\mathbf{x}_k + \alpha_k \mathbf{d}_k) \leq f(\mathbf{x}_k) + \rho \alpha_k \mathbf{g}_k^T \mathbf{d}_k$$

$$\nabla f(\mathbf{x}_k + \alpha_k \mathbf{d}_k)^T \mathbf{d}_k \geq \sigma \mathbf{g}_k^T \mathbf{d}_k$$

其中: $\mathbf{g}_k = \mathbf{g}(\mathbf{x}_k) = \nabla f(\mathbf{x}_k)$

有时用更强的条件代替: (强Wolfe准则)

$$|\nabla f(\mathbf{x}_k + \alpha_k \mathbf{d}_k)^T \mathbf{d}_k| \leq -\sigma \mathbf{g}_k^T \mathbf{d}_k$$

$\sigma > 0$ 充分小时, 可使其变为近似精确线搜索.

由该Wolfe准则得到的新的迭代点 $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k$ 在 \mathbf{x}_k 的某一领域内并且使目标函数值有一定的下降量.

由于 $\mathbf{g}_k^T \mathbf{d}_k < 0$, 有定理可以证明Wolfe准则的有限终止性.

非精确线搜索方法——Armijo 准则

- Armijo 准则是指给定 $\beta \in (0, 1)$, $\sigma \in (0, 0.5)$. 令步长因子 $\alpha_k = \beta^{m_k}$, 其中 m_k 为满足下列不等式的最小非负整数:

$$f(\mathbf{x}_k + \beta^m \mathbf{d}_k) \leq f(\mathbf{x}_k) + \sigma \beta^m \mathbf{g}_k^T \mathbf{d}_k$$

可以证明, 若 $f(\mathbf{x})$ 是连续可微的且满足 $\mathbf{g}_k^T \mathbf{d}_k < 0$, 则 Armijo 准则是有限终止的, 即存在正数 σ , 使得对于充分大的正整数 m , 上式成立.

算法 2.4 (Armijo 准则)

给定参数 $\beta \in (0, 1)$, $\sigma \in (0, 0.5)$. 对 $m = 0, 1, \dots$, 若不等式

$$f(\mathbf{x}_k + \beta^m \mathbf{d}_k) \leq f(\mathbf{x}_k) + \sigma \beta^m \mathbf{g}_k^T \mathbf{d}_k$$

成立, 则置 $m_k := m$, $\mathbf{x}_{k+1} := \mathbf{x}_k + \beta^{m_k} \mathbf{d}_k$, 第 k 步线搜索完成.

非精确线搜索方法——Armijo准则

● Armijo搜索的Matlab程序:

```
function mk=armijo(xk,dk)
beta=0.5; sigma=0.2;
m=0; maxm=20;
while (m<=maxm)
    if (fun(xk+beta^m*dk) <= fun(xk) + sigma*beta^m*gfun(xk)'*dk)
        mk=m; break;
    end
    m=m+1;
end
alpha=beta^mk
newxk=xk+alpha*dk
fk=fun(xk)
newfk=fun(newxk)
```

fun: 目标函数

gfun: 目标函数的梯度函数

线搜索法总结

- **线搜索法**: 指用**线搜索策略求步长因子**的无约束优化问题下降类算法的简称, 一般算法框架如下.

算法 2.5 (线搜索法算法框架)

步骤 0, 初始化. 选取有关参数及初始迭代点 $\mathbf{x}_0 \in \mathbf{R}^n$. 设定容许误差 $0 \leq \varepsilon \ll 1$. 令 $k := 0$.

步骤 1, 检验终止判别准则. 计算 $\mathbf{g}_k = \nabla f(\mathbf{x}_k)$. 若 $\|\mathbf{g}_k\| \leq \varepsilon$, 停算, 输出 $\mathbf{x}^* \approx \mathbf{x}_k$.

步骤 2, 确定下降方向 \mathbf{d}_k , 使满足 $\mathbf{g}_k^T \mathbf{d}_k < 0$.

步骤 3, 确定步长因子 α_k . 可在下列“精确”与“非精确”两种线搜索策略中选用其一:

(1) 用黄金分割法或抛物线法等精确线搜索策略求

$$\alpha_k = \arg \min_{\alpha > 0} f(\mathbf{x}_k + \alpha \mathbf{d}_k);$$

(2) 用 Wolfe 准则或 Armijo 准则等非精确线搜索策略求 α_k

步骤 4, 更新迭代点. 令 $\mathbf{x}_{k+1} := \mathbf{x}_k + \alpha_k \mathbf{d}_k$, $k := k + 1$, 转步骤 1.

梯度法

- 梯度法(最速下降法)是求解无约束优化问题最简单最古老的方法之一. 前述下降类算法的一般框架中提到用不同的方式确定搜索方向或搜索步长, 就会得到不同算法. 梯度法是用负梯度方向

$$\mathbf{d}_k = -\nabla f(\mathbf{x}_k)$$

作为搜索方向的.

- 两个问题:
 - ✓ 为什么选择负梯度方向?
 - ✓ 为什么称为最速下降法?

梯度法

1) 为什么选择负梯度方向? 2) 为什么称为最速下降法?

搜索步长

设 $f(\mathbf{x})$ 在 \mathbf{x}_k 附近连续可微, \mathbf{d}_k 为搜索方向向量, $\mathbf{g}_k = \nabla f(\mathbf{x}_k)$. 泰勒展开得

$$f(\mathbf{x}_k + \alpha \mathbf{d}_k) = f(\mathbf{x}_k) + \alpha \mathbf{g}_k^T \mathbf{d}_k + o(\alpha), \quad \alpha > 0.$$

那么目标函数 $f(\mathbf{x})$ 在 \mathbf{x}_k 处沿方向 \mathbf{d}_k 下降的变化率为

$$\begin{aligned} \lim_{\alpha \rightarrow 0} \frac{f(\mathbf{x}_k + \alpha \mathbf{d}_k) - f(\mathbf{x}_k)}{\alpha} &= \lim_{\alpha \rightarrow 0} \frac{\alpha \mathbf{g}_k^T \mathbf{d}_k + o(\alpha)}{\alpha} \\ &= \mathbf{g}_k^T \mathbf{d}_k = \|\mathbf{g}_k\| \|\mathbf{d}_k\| \cos \bar{\theta}_k, \end{aligned}$$

$\bar{\theta}_k$ 为 \mathbf{g}_k 与 \mathbf{d}_k 的夹角.

小于0

也即搜索方向为负梯度方向, 是目标函数在当前点的最速下降方向

要使变化率最小, 也即变化率绝对值最大, 只有 $\cos \bar{\theta}_k = -1$, 即 $\bar{\theta}_k = \pi$

梯度法

● 梯度法计算步骤:

用精确线搜索方法

$$\phi'(\alpha) = \frac{d}{d\alpha} f(\mathbf{x}_k + \alpha \mathbf{d}_k) \Big|_{\alpha=\alpha_k} = \nabla f(\mathbf{x}_k + \alpha_k \mathbf{d}_k)^\top \mathbf{d}_k = 0.$$

有
$$\mathbf{g}(\mathbf{x}_{k+1})^\top \mathbf{g}(\mathbf{x}_k) = 0,$$

即新点 \mathbf{x}_{k+1} 处的梯度与旧点 \mathbf{x}_k 处的梯度是正交的

算法 3.1 (梯度法)

步骤 0, 选取初始点 $\mathbf{x}_0 \in \mathbf{R}^n$, 容许误差 $0 \leq \varepsilon \ll 1$. 令 $k := 0$.

步骤 1, 计算 $\mathbf{g}_k = \nabla f(\mathbf{x}_k)$. 若 $\|\mathbf{g}_k\| \leq \varepsilon$, 停算, 输出 \mathbf{x}_k 作为近似极小点.

步骤 2, 取方向 $\mathbf{d}_k = -\mathbf{g}_k$.

步骤 3, 由线搜索方法确定步长因子 α_k .

步骤 4, 令 $\mathbf{x}_{k+1} := \mathbf{x}_k + \alpha_k \mathbf{d}_k$, $k := k + 1$, 转步骤 1.

可用精确/非精确线性搜索方法, 在理论上均能保证全局收敛性.

梯度法

● 一些定理(了解)

定理 3.1 设目标函数 $f(\mathbf{x})$ 连续可微且其梯度函数 $\mathbf{g}(\mathbf{x})$ 是 Lipschitz 连续的, $\{\mathbf{x}_k\}$ 由梯度法产生, 其中步长因子 α_k 由精确线搜索, 或由 Wolfe 准则, 或由 Armijo 准则产生, 则有

$$\lim_{k \rightarrow \infty} \|\mathbf{g}(\mathbf{x}_k)\| = 0.$$

定理 3.2 设矩阵 $\mathbf{A} \in \mathbf{R}^{n \times n}$ 对称正定, $\mathbf{b} \in \mathbf{R}^n$. 记 λ_1 和 λ_n 分别是 \mathbf{A} 的最大和最小特征值, $\kappa = \lambda_1/\lambda_n$. 考虑如下极小化问题

$$\min f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} - \mathbf{b}^T \mathbf{x}.$$

设 $\{\mathbf{x}_k\}$ 是用精确线搜索的梯度法求解上述问题所产生的迭代序列, 则对于所有的 k , 下面的不等式成立

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\|_A \leq \left(\frac{\kappa - 1}{\kappa + 1} \right) \|\mathbf{x}_k - \mathbf{x}^*\|_A,$$

式中: \mathbf{x}^* 为问题的唯一解, $\|\mathbf{x}\|_A = \sqrt{\mathbf{x}^T \mathbf{A} \mathbf{x}}$.

κ 接近1(\mathbf{A} 的最大和最小特征值接近时), 梯度法收敛很快. κ 较大时(\mathbf{A} 近似于病态), 收敛很慢

牛顿法

- **牛顿法基本思想**: 用迭代点 x_k 处的一阶导数(梯度)和二阶导数(Hesse矩阵)对目标函数进行二次函数近似, 然后把二次函数的极小点作为新的迭代点, 不断重复这一过程, 直至求得满足精度的近似极小点.

设 $f(x)$ 的 Hesse 阵 $G(x) = \nabla^2 f(x)$ 连续, 取其在 x_k 处的泰勒展开式的前三项, 得

$$q_k(x) = f_k + \mathbf{g}_k^T(x - x_k) + \frac{1}{2}(x - x_k)^T \mathbf{G}_k(x - x_k),$$

求二次函数 $q_k(x)$ 的稳定点, 得

$$\nabla q_k(x) = \mathbf{g}_k + \mathbf{G}_k(x - x_k) = \mathbf{0}.$$

若 G_k 非奇异, 那么解上面的线性方程组 得

$$x_{k+1} = x_k - G_k^{-1} g_k.$$

实际运算中为避免求逆, 可先解 $G_k d = -g_k$ 得 $d_k = x_{k+1} - x_k$, 然后令 $x_{k+1} = x_k + d_k$

牛顿法

- 基本牛顿法**优点**: 收敛速度快, 具有局部二阶收敛性

算法 3.2 (基本牛顿法)

步骤 0, 选取初始点 $\mathbf{x}_0 \in \mathbf{R}^n$, 容许误差 $0 \leq \varepsilon \ll 1$. 令 $k := 0$.

步骤 1, 计算 $\mathbf{g}_k = \nabla f(\mathbf{x}_k)$. 若 $\|\mathbf{g}_k\| \leq \varepsilon$, 停算, 输出 $\mathbf{x}^* \approx \mathbf{x}_k$.

步骤 2, 计算 $\mathbf{G}_k = \nabla^2 f(\mathbf{x}_k)$, 并求解线性方程组

$$\mathbf{G}_k \mathbf{d} = -\mathbf{g}_k,$$

得解 \mathbf{d}_k .

步骤 3, 令 $\mathbf{x}_{k+1} := \mathbf{x}_k + \mathbf{d}_k$, $k := k + 1$, 转步骤 1.

牛顿法

- **基本牛顿法缺点**: 初始点需要足够靠近极小点, 否则算法有可能不收敛. 但通常精确极小点位置未知, 故而带来困难. 可引入线搜索方法以得到大范围收敛的算法 → **阻尼牛顿法**

定理 3.3 设函数 $f(\mathbf{x})$ 有二阶连续偏导数, 在局部极小点 \mathbf{x}^* 处, $\mathbf{G}(\mathbf{x}^*) = \nabla^2 f(\mathbf{x}^*)$ 是正定的且 $\mathbf{G}(\mathbf{x})$ 在 \mathbf{x}^* 的一个邻域内是 Lipschitz 连续的. 如果 **初始点 \mathbf{x}_0 充分靠近 \mathbf{x}^*** , 那么对一切 k , 牛顿迭代公式 **(3.4)** 是适定的, 并且当 $\{\mathbf{x}_k\}$ 为无穷点列时, 其极限为 \mathbf{x}^* 且收敛阶 **至少是二阶的.**

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{G}_k^{-1} \mathbf{g}_k.$$

牛顿法

● 基于Armijo搜索的阻尼牛顿法:

算法 3.3 (基于 Armijo 搜索的阻尼牛顿法)

步骤 0, 选取参数 $\beta \in (0, 1)$, $\sigma \in (0, 0.5)$, 初始点 $\mathbf{x}_0 \in \mathbf{R}^n$, 容许误差 $0 \leq \varepsilon \ll 1$. 令 $k := 0$.

步骤 1, 计算 $\mathbf{g}_k = \nabla f(\mathbf{x}_k)$. 若 $\|\mathbf{g}_k\| \leq \varepsilon$, 停算, 输出 $\mathbf{x}^* \approx \mathbf{x}_k$.

步骤 2, 计算 $\mathbf{G}_k = \nabla^2 f(\mathbf{x}_k)$, 并求解线性方程组

$$\mathbf{G}_k \mathbf{d} = -\mathbf{g}_k,$$

牛顿法求
搜索方向

得解 \mathbf{d}_k .

步骤 3, 记 m_k 是满足下列不等式的最小非负整数 m :

$$f(\mathbf{x}_k + \beta^m \mathbf{d}_k) \leq f(\mathbf{x}_k) + \sigma \beta^m \mathbf{g}_k^T \mathbf{d}_k.$$

Armijo 法则
搜索步长

步骤 4, 令 $\alpha_k := \beta^{m_k}$, $\mathbf{x}_{k+1} := \mathbf{x}_k + \alpha_k \mathbf{d}_k$, $k := k + 1$, 转步骤 1.

牛顿法

● 基于Armijo搜索的阻尼牛顿法(理论支撑,了解)

定理 3.4 设函数 $f(\mathbf{x})$ 二次连续可微且存在常数 $\gamma > 0$, 使得

$$\mathbf{d}^T \mathbf{G}(\mathbf{x}) \mathbf{d} \geq \gamma \|\mathbf{d}\|^2, \quad \forall \mathbf{d} \in \mathbf{R}^n, \mathbf{x} \in \mathcal{L}(\mathbf{x}_0), \quad (3.8)$$

式中: $\mathcal{L}(\mathbf{x}_0) = \{\mathbf{x} | f(\mathbf{x}) \leq f(\mathbf{x}_0)\}$. 设 $\{\mathbf{x}_k\}$ 是由算法 3.3 产生的无穷点列, 则该点列收敛于 f 在水平集 $\mathcal{L}(\mathbf{x}_0)$ 中的唯一全局极小点.

引理 3.1 设函数 $f: \mathbf{R}^n \rightarrow \mathbf{R}$ 二次连续可微, 点列 $\{\mathbf{x}_k\}$ 由算法 3.3 产生. 设 $\{\mathbf{x}_k\} \rightarrow \mathbf{x}^*$ 且 $\mathbf{g}(\mathbf{x}^*) = \mathbf{0}$, $\mathbf{G}(\mathbf{x}^*)$ 正定. 那么, 若

$$\lim_{k \rightarrow \infty} \frac{\|\mathbf{G}(\mathbf{x}_k) \mathbf{d}_k + \mathbf{g}_k\|}{\|\mathbf{d}_k\|} = 0, \quad (3.10)$$

- 则
- (1) 当 k 充分大时, 步长因子 $\alpha_k \equiv 1$;
 - (2) 点列 $\{\mathbf{x}_k\}$ 超线性收敛于 \mathbf{x}^* .

定理 3.5 设定理 3.4 的条件成立, 点列 $\{\mathbf{x}_k\}$ 由算法 3.3 产生, 则 $\{\mathbf{x}_k\}$ 超线性收敛于 f 的全局极小点 \mathbf{x}^* . 此外, 若 Hesse 阵 $\mathbf{G}(\mathbf{x})$ 在 \mathbf{x}^* 处 Lipschitz 连续, 则收敛阶至少是二阶的.

修正牛顿法

- 牛顿法的**优点**: 不低于二阶的收敛速度
- 牛顿法的**缺点**: 要求目标函数的Hesse矩阵 $\mathbf{G}(\mathbf{x})$ 在每个迭代点 \mathbf{x}_k 处是**正定**的, 否则难以保证牛顿方向 $\mathbf{d}_k = -\mathbf{G}_k^{-1} \mathbf{g}_k$ 是 f 在 \mathbf{x}_k 处的下降方向.
- **修正方法之一**: 和梯度法结合起来, 构造“**牛顿-梯度混合算法**”: 当Hesse阵正定, 用牛顿方向作为搜索方向; 否则, 若Hesse阵奇异, 或非奇异但牛顿方向不是下降方向, 采用负梯度方向作为搜索方向.

修正牛顿法--牛顿-梯度混合算法

算法 3.4 (牛顿-梯度混合算法)

步骤 0, 选取初始点 $\mathbf{x}_0 \in \mathbf{R}^n$, 容许误差 $0 \leq \varepsilon \ll 1$. 令 $k := 0$.

步骤 1, 计算 $\mathbf{g}_k = \nabla f(\mathbf{x}_k)$. 若 $\|\mathbf{g}_k\| \leq \varepsilon$, 停算, 输出 \mathbf{x}_k 作为近似极小点.

步骤 2, 计算 $\mathbf{G}_k = \nabla^2 f(\mathbf{x}_k)$, 并求解线性方程组

$$\mathbf{G}_k \mathbf{d} = -\mathbf{g}_k.$$

若方程组

有解 \mathbf{d}_k 且满足 $\mathbf{g}_k^T \mathbf{d}_k < 0$, 转步骤 3;

否则, 令 $\mathbf{d}_k = -\mathbf{g}_k$, 转步骤 3.

步骤 3, 由线搜索方法确定步长因子 α_k .

步骤 4, 令 $\mathbf{x}_{k+1} := \mathbf{x}_k + \alpha_k \mathbf{d}_k$, $k := k + 1$, 转步骤 1.

修正牛顿法

- 牛顿法的**优点**: 不低于二阶的收敛速度
- 牛顿法的**缺点**: 要求目标函数的Hesse矩阵 $\mathbf{G}(\mathbf{x})$ 在每个迭代点 \mathbf{x}_k 处是**正定**的, 否则难以保证牛顿方向 $\mathbf{d}_k = -\mathbf{G}_k^{-1} \mathbf{g}_k$ 是 f 在 \mathbf{x}_k 处的下降方向.
- **修正方法之二**: 引入阻尼因子 $\mu_k \geq 0$, 即在每一迭代步适当选取参数 μ_k 使得矩阵 $\mathbf{A}_k = \mathbf{G}(\mathbf{x}_k) + \mu_k \mathbf{I}$ 正定.

修正牛顿法

算法 3.5 (修正牛顿法)

步骤 0, 选取参数 $\beta \in (0, 1)$, $\sigma \in (0, 0.5)$, 初始点 $\mathbf{x}_0 \in \mathbf{R}^n$, 容许误差 $0 \leq \varepsilon \ll 1$, 参数 $\tau \in [0, 1]$. 令 $k := 0$.

步骤 1, 计算 $\mathbf{g}_k = \nabla f(\mathbf{x}_k)$, $\mu_k = \|\mathbf{g}_k\|^{1+\tau}$. 若 $\|\mathbf{g}_k\| \leq \varepsilon$, 停算, 输出 \mathbf{x}_k 作为近似极小点.

步骤 2, 计算 Hesse 阵 $\mathbf{G}_k = \nabla^2 f(\mathbf{x}_k)$, 并求解线性方程组

$$(\mathbf{G}_k + \mu_k \mathbf{I})\mathbf{d} = -\mathbf{g}_k$$

牛顿法求
搜索方向

得解 \mathbf{d}_k .

步骤 3, 令 m_k 是满足下列不等式的最小非负整数 m :

$$f(\mathbf{x}_k + \beta^m \mathbf{d}_k) \leq f(\mathbf{x}_k) + \sigma \beta^m \mathbf{g}_k^T \mathbf{d}_k.$$

Armijo 法则
搜索步长

步骤 4, 令 $\alpha_k := \beta^{m_k}$, $\mathbf{x}_{k+1} := \mathbf{x}_k + \alpha_k \mathbf{d}_k$, $k := k + 1$, 转步骤 1.