

梯度下降法 (最速下降法)

$$x^{k+1} = x^k - \alpha \nabla f(x^k)$$

Then: ϕ 二次正定的函数 $f(x) = \frac{1}{2} x^T A x - b^T x$

$\phi \rightarrow$ 精确搜索步 $\arg \min_{\alpha > 0} \phi(\alpha) = f(x^k - \alpha \nabla f(x^k))$

$$\alpha_k = \frac{\|\nabla f(x^k)\|^2}{\nabla f(x^k)^T A \nabla f(x^k)}$$

② Q-线性收敛: $\|x^{k+1} - x^*\|_A \leq \left(\frac{\lambda_n - \lambda_1}{\lambda_n + \lambda_1} \right) \|x^k - x^*\|_A$

其中 λ_1, λ_n 是 A 的最小, 最大特征值, $\|x\|_A = \sqrt{x^T A x}$

Proof:

$$x^{k+1} - x^* = x^k - x^* - \alpha \nabla f(x^k)$$

$$= x^k - x^* - \alpha (A x^k - b)$$

$$= x^k - x^* - \alpha A (x^k - x^*)$$

$$\nabla f(x^k) = A x^k - b$$

$$\nabla f(x^*) = 0$$

$$= A x^* - b = 0$$

$$e^k = x^k - x^*$$

$$e^{k+1} = (I - \alpha A) e^k$$

$$\|x\|_A^2 = x^T A x$$

$$\|e^{k+1}\|_A^2 = (e^{k+1})^T A e^{k+1}$$

$$= e^{kT} (I - \alpha A)^T A (I - \alpha A) e^k$$

$$= e^{kT} (A - 2\alpha A^2 + \alpha^2 A^3) e^k$$

$$= e^{kT} Q (I - 2\alpha \Lambda^2 + \alpha^2 \Lambda^3) Q^T e^k$$

$$= \sum_{i=1}^n \lambda_i (1 - 2\alpha \lambda_i + \alpha^2 \lambda_i^2) y_i^2$$

$$A = Q \Lambda Q^T \quad Q^T Q = I$$

$$A^2 = Q \Lambda Q^T Q \Lambda Q^T$$

$$= Q \Lambda^2 Q^T$$

$$A^3 = Q \Lambda^3 Q^T$$

$$y = Q^T e^k$$

$$\sum_{i=1}^n \lambda_i y_i^2$$

$$= \sum_{i=1}^n \lambda_i (Q^T e^k)_i^2$$

$$= e^{kT} A e^k$$

$$= \|e^k\|_A^2$$

由于

$$\|A x^k - b\|_A^2 = \text{tr} \left((A x^k - b)^T A (A x^k - b) \right) \\ = \text{tr} \left((A x^k - b)^T Q \Lambda Q^T (A x^k - b) \right)$$

$$\lambda_1 \|A x^k - b\|_2 \leq \|A x^k - b\|_A \leq \lambda_n \|A x^k - b\|_2$$

$$\text{又因为 } \alpha = \frac{\|\nabla f(x^k)\|_2^2}{\nabla f(x^k)^T A \nabla f(x^k)} = \frac{\|A x^k - b\|_2^2}{\|A x^k - b\|_A^2} \geq \frac{2}{\lambda_1 + \lambda_n}$$

$$1 - \alpha \lambda_i \leq 1 - \frac{2}{\lambda_1 + \lambda_n} \lambda_i \leq 1 - \frac{2}{\lambda_1 + \lambda_n} \lambda_n = \frac{\lambda_1 - \lambda_n}{\lambda_1 + \lambda_n}$$

$$\|e^{k+1}\|_A^2 \leq \left(\frac{\lambda_1 - \lambda_n}{\lambda_1 + \lambda_n} \right)^2 \|e^k\|_A^2$$

Defn (Condition Number 条件数) 对于可逆矩阵 A , 其条件数为

$$\kappa(A) = \|A\| \cdot \|A^{-1}\| \rightarrow \text{范数}$$

最常用: 基于谱范数 (spectral norm)

$$\kappa(A) = \frac{\sigma_{\max}(A)}{\sigma_{\min}(A)}$$

对于正定矩阵: $\kappa(A) = \frac{\lambda_{\max}}{\lambda_{\min}}$

$\kappa(A) \downarrow$, 收敛 \downarrow $\kappa(A) \downarrow$, 收敛 \uparrow

最优性条件

梯度下降法

初始化: 初始点 x^0 , 步长 α , 最大迭代数 K . $\varepsilon \in (0, 1)$

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

$$x^{k+1} = x^k - \alpha \nabla f(x^k)$$

If $\|\nabla f(x^k)\| < \varepsilon$ then

STOP

近似一阶最优性条件.

End If

End For

输出 x^{k+1}

Thm (一阶必要条件) 假设 f 在全空间 \mathbb{R}^n 可微, 如果 x^* 是一个局部极小点, 那么

$$\nabla f(x^*) = 0$$

Proof: 设 x^* 为局部极小点,

$$\forall d \in \mathbb{R}^n \quad f(x^* + td) = f(x^*) + \nabla f(x^*)^T(td) + o(t)$$

$$\Rightarrow \frac{f(x^* + td) - f(x^*)}{t} = d^T \nabla f(x^*) + o(1)$$

$$\lim_{t \rightarrow 0^+} \frac{f(x^* + td) - f(x^*)}{t} = d^T \nabla f(x^*) \geq 0 \quad \forall d \in \mathbb{R}^n$$

$$\lim_{t \rightarrow 0^-} \frac{f(x^* + td) - f(x^*)}{t} = d^T \nabla f(x^*) \leq 0$$

$$\Rightarrow d^T \nabla f(x^*) = 0 \Rightarrow \nabla f(x^*) = 0$$

→ 也称为稳定点 (stationary point)

① 一阶必要条件.

$$f(x) = x^2, \quad f(x) = -x^2$$

$$f(x) = x^3$$

$$f(x, y) = x^2 - y^2$$

Thm (二阶最优性条件) 假设 f 在点 x^* 的一个开邻域内是二阶连续可微的. 则:

① (一阶最优性) 若 x^* 是 f 的一个局部极小点, 则

$$\nabla f(x^*) = 0, \quad \nabla^2 f(x^*) \succeq 0$$

$$[\nabla^2 f(x^*)]_{ij} = \frac{\partial^2 f(x^*)}{\partial x_i \partial x_j}$$

② (充分性) 若存在 x^* , 有

$$\nabla f(x^*) = 0, \quad \nabla^2 f(x^*) \succ 0$$

则 x^* 是一个局部极小点.

← 全局最优性.

Proof: 证 ①, 设 x^* 为局部极小点, 在 x^* 处二阶泰勒展开:

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

由于 $\nabla f(x^*) = 0$. 反设 $\nabla^2 f(x^*) \not\succeq 0$ 不成立. 即存在特征值 $\lambda_- < 0$.

令 d 为 λ_- 对应的特征向量.

$$\frac{f(x^* + d) - f(x^*)}{\|d\|^2} = \frac{1}{2} \frac{1}{\|d\|^2} d^T \nabla^2 f(x^*) d + o(1) = \frac{1}{2} \lambda_- + o(1)$$

$$\nabla^2 f(x^*) d = \lambda_- d$$

$$d^T \nabla^2 f(x^*) d = d^T \lambda_- d$$

$$= \lambda_- d^T d$$

$$= \lambda_- \|d\|^2$$

$$\text{当 } \|d\| \rightarrow 0, \quad \frac{f(x^* + d) - f(x^*)}{\|d\|^2} = \frac{1}{2} \lambda_- < 0$$

$$\Rightarrow f(x^* + d) < f(x^*) \quad x^* \text{ 最优, 矛盾}$$

② (充分性) 若存在 x^* , 有

$$\nabla f(x^*) = 0, \quad \nabla^2 f(x^*) \succ 0$$

则 x^* 是 f 局部极小点. ← 全局最优性.

$$\text{设 } \nabla^2 f(x^*) \succ 0, \quad d^T \nabla^2 f(x^*) d \geq \lambda_{\min} \|d\|^2 > 0 \quad \forall d \neq 0$$

$$\Rightarrow \frac{f(x^* + d) - f(x^*)}{\|d\|^2} = \frac{1}{2} \frac{1}{\|d\|^2} d^T \nabla^2 f(x^*) d \geq \frac{1}{2} \lambda_{\min} + o(1)$$

当 $\|d\| \rightarrow 0$

$$\frac{f(x^* + d) - f(x^*)}{\|d\|^2} \geq \frac{1}{2} \lambda_{\min} > 0$$

$$\Rightarrow f(x^* + d) > f(x^*) \quad \forall d \neq 0 \quad \therefore f(x^*) \text{ 是局部最优}$$



Defn (梯度 Lipschitz 连续) 设 f 为可微函数, 若存在 $L > 0$,

$$\|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\| \quad \forall x, y \in \text{dom} f$$

则称 f 是 (全局的) L -梯度 Lipschitz 连续函数

$$\text{E.g. } f(x) = x^2 \quad \|\nabla f(x) - \nabla f(y)\| = \|2x - 2y\| \Rightarrow \|x - y\|$$

$$f(x) = x^3 \quad \|\nabla f(x) - \nabla f(y)\| = \|3x^2 - 3y^2\| \leq 3\|x+y\| \|x-y\|$$

→ 局部梯度 Lipschitz 连续, 非全局...

$$f(x) = e^x \quad x \in [0, 100] \quad \text{微分中值定理}$$

$$\|\nabla f(x) - \nabla f(y)\| = \|\nabla^2 f(\xi)\| \|x - y\| \quad \xi \in [0, 100]$$

$$\leq \|\nabla^2 f(\xi)\| \|x - y\|$$

$$L = \max_{\xi \in [0, 100]} \|\nabla^2 f(\xi)\| = e^{100}$$

梯度 Lipschitz 表明 $\nabla f(x)$ 的变化可被 x 的变化控制.

Lemma (二次上界) 设 f 为可微函数, $\text{dom } f \subset \mathbb{R}^n$, 且 f 是 L -梯度 Lipschitz 连续函数, 则

$$f(y) \leq f(x) + \nabla f(x)^T (y-x) + \frac{1}{2} L \|y-x\|^2 \quad \forall x, y \in \text{dom } f$$

\Rightarrow 二次上界.

① 要求 f 的增量不超过二次.

Then (梯度法在凸函数上的收敛性) 设 f 为凸, 梯度 L -Lipschitz 连续, $f^* = \inf_x f(x)$ 可达. 如果步长 α 为常数且 $0 < \alpha < \frac{1}{L}$, 那么.

- ① 梯度下降法生成的点列 $\{x^k\}$ 收敛到最优值 \rightarrow 全局最优
- ② 函数值表义下收敛速度为 $O(\frac{1}{k})$

Proof: 令 $\tilde{x} = x - \alpha \nabla f(x)$

$$\begin{aligned} f(\tilde{x}) &\leq f(x) + \nabla f(x)^T (-\alpha \nabla f(x)) + \frac{1}{2} \|\alpha \nabla f(x)\|^2 \\ &= f(x) - \alpha \left(1 - \frac{L\alpha}{2}\right) \|\nabla f(x)\|^2 \end{aligned}$$

$$\text{令 } 0 < \alpha < \frac{1}{L} \Rightarrow \left(1 - \frac{L\alpha}{2}\right) > \frac{1}{2}$$

$$\therefore f(\tilde{x}) \leq f(x) - \alpha \left(1 - \frac{L\alpha}{2}\right) \|\nabla f(x)\|^2$$

$$\begin{aligned} f^* &\leq f(x) + \nabla f(x)^T (x^* - x) \\ &\Rightarrow f(x) \leq f^* + \nabla f(x)^T (x - x^*) \end{aligned}$$

$$\leq f(x) - \frac{\alpha}{2} \|\nabla f(x)\|^2$$

$$\leq f^* + \nabla f(x)^T (x - x^*) - \frac{\alpha}{2} \|\nabla f(x)\|^2$$

$$= f^* + \frac{1}{2\alpha} \left(\|x - x^*\|^2 - \|x - x^*\|^2 + 2\alpha \nabla f(x)^T (x - x^*) - \alpha^2 \|\nabla f(x)\|^2 \right)$$

$$= f^* + \frac{1}{2\alpha} \left(\|x - x^*\|^2 - \|x - x^* - \alpha \nabla f(x)\|^2 \right)$$

$$= f^* + \frac{1}{2\alpha} \left(\|x - x^*\|^2 - \|\tilde{x} - x^*\|^2 \right)$$

$$\Rightarrow f(\tilde{x}) - f^* \leq \frac{1}{2\alpha} \left(\|x - x^*\|^2 - \|\tilde{x} - x^*\|^2 \right)$$

$x = x^k, \tilde{x} = x^{k+1}, k=0, \dots, K$, 相加:

$$\Rightarrow \sum_{k=0}^K (f(x^{k+1}) - f^*) \leq \frac{1}{2\alpha} \sum_{k=0}^K \left(\|x^k - x^*\|^2 - \|x^{k+1} - x^*\|^2 \right)$$

$$= \frac{1}{2\alpha} \left(\|x^0 - x^*\|^2 - \|x^{K+1} - x^*\|^2 \right)$$

$$\leq \frac{1}{2\alpha} \|x^0 - x^*\|^2$$

$$\Rightarrow f(x^{K+1}) - f^* \leq \frac{1}{K} \sum_{k=1}^K (f(x^k) - f^*) \leq \frac{1}{K} \frac{1}{2\alpha} \|x^0 - x^*\|^2$$