

运筹学与优化方法

晁国清

计算机科学与技术学院



课程回顾

➤ 无约束优化问题的最优性条件

若 $f(x)$ 为一般函数且 $f(x)$ 在 x^* 处一阶可微, 则

$$x^* \text{ 是最优解} \Rightarrow \nabla f(x^*) = 0$$

若 $f(x)$ 为一般函数且 $f(x)$ 在 x^* 处二阶可微, 则

$$x^* \text{ 是最优解} \Rightarrow \nabla^2 f(x^*) \geq 0$$

若 $f(x)$ 为一般函数且 $f(x)$ 在 x^* 处二阶可微

$$\nabla f(x^*) = 0 \text{ 且 } \nabla^2 f(x^*) > 0 \Rightarrow x^* \text{ 是最优解}$$

➤ 迭代算法的基本思想

给定初始点 x^0 , 产生点列

$$\{x^k\}_{k=1}^{\infty},$$

并且点列满足条件

$$f(x^{k+1}) < f(x^k)$$

Step 1: 给定初始点 x_0 , $k=0$;

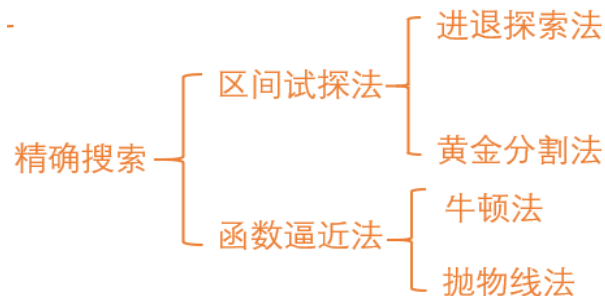
Step 2: 判断 x^k 是否满足终止条件; 是, 结束; 不是, 进行Step 3;

Step 3: 寻找 x^k 处的下降方向 d^k ;

Step 4: 选择合适的步长 $\alpha_k > 0$, 使 $f(x^k + \alpha_k d^k) < f(x^k)$ 成立;

Step 5: 令 $x^{k+1} = x^k + \alpha_k d^k$ 且 $k=k+1$; 转Step 2

➤ 一维(线性)搜索



课程内容

- 非精确的线搜索
- 变量轮换法
- 最速下降法



非精确线搜索

➤Armijo-Goldstein准则：定义 $\phi(\alpha) = f(\mathbf{x}_k + \alpha \mathbf{d}_k)$, 则有

$$\phi(0) + (1 - \rho)\alpha\phi'(0) \leq \phi(\alpha) \leq \phi(0) + \rho\alpha\phi'(0), \rho \in (0, 0.5) \quad (*)$$

假设 \mathbf{d}_k 是下降方向, 对 $f(\mathbf{x}_k + \alpha \mathbf{d}_k)$ 在 \mathbf{x}_k 处一阶泰勒展开:

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

已知 \mathbf{d}_k 是下降方向, 则有 $\nabla f(\mathbf{x}_k)^T \mathbf{d}_k < 0$, 为保证下降, 找 $f(\mathbf{x}_k) - f(\mathbf{x}_k + \alpha \mathbf{d}_k)$ 的一个合理下界:

$$f(\mathbf{x}_k) - f(\mathbf{x}_k + \alpha \mathbf{d}_k) \geq -\rho\alpha \nabla f(\mathbf{x}_k)^T \mathbf{d}_k, \rho \in (0, 0.5) \quad (1)$$

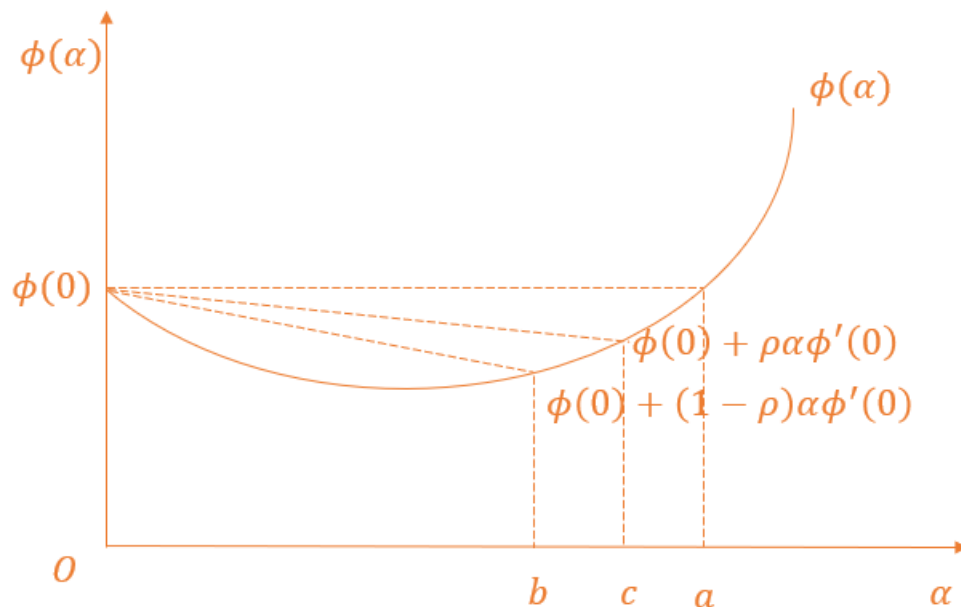
再给其一个上界(α 不能太小)

$$f(\mathbf{x}_k) - f(\mathbf{x}_k + \alpha \mathbf{d}_k) \leq -(1 - \rho)\alpha \nabla f(\mathbf{x}_k)^T \mathbf{d}_k, \rho \in (0, 0.5) \quad (2)$$

将(1)和(2)写在一起, 并定义 $\phi(\alpha) = f(\mathbf{x}_k + \alpha \mathbf{d}_k)$, 就得到(*)



Armijo-Goldstein准则的几何意义



从此图可以看到， $\phi(0) + \rho\alpha\phi'(0)$ 在 $\phi(0)$ 下方，在 $\phi(0) + (1 - \rho)\alpha\phi'(0)$ 的上方，对应区间 $[b, c]$ 就是满足Armijo-Goldstein准则的步长

缺点： Armijo-Goldstein准则可能会把极值点排除在可接受的区间之外

Wolfe-Powell准则

- Wolfe-Powell准则是为了解决Armijo-Goldstein准则可能会把极值点排除在可接受的区间之外的问题提出来的
- Wolfe-Powell准则的上界与Armijo-Goldstein准则是一样的，即
$$f(\mathbf{x}_k + \alpha \mathbf{d}_k) - f(\mathbf{x}_k) \leq \rho \alpha \nabla f(\mathbf{x}_k)^T \mathbf{d}_k, \rho \in (0, 0.5)$$

为了保证步长足够大，并且可接受的区间包含极值点，上界被定义为：

$$\nabla f(\mathbf{x}_{k+1})^T \mathbf{d}_k \geq \sigma \nabla f(\mathbf{x}_k)^T \mathbf{d}_k, \sigma \in (\rho, 1)$$

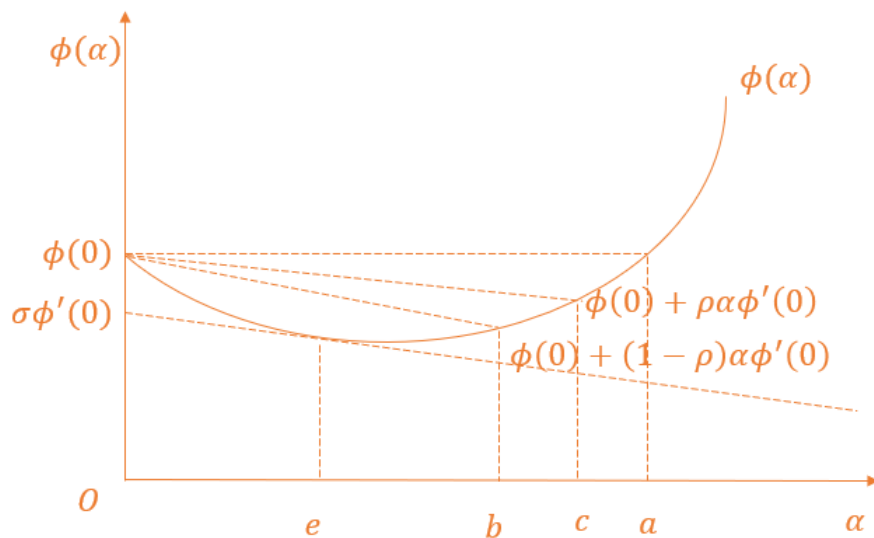
即 $\phi'(\alpha) \geq \sigma \phi'(0)$ ，说明在 α 处的斜率应该大于起始点斜率的 σ 倍。 $\phi'(0)$ 是负值，所以上界的含义就是可接受范围中斜率接近于零的负值或者正值。

强Wolfe条件： $|\phi'(\alpha)| \leq -\sigma \phi'(0)$

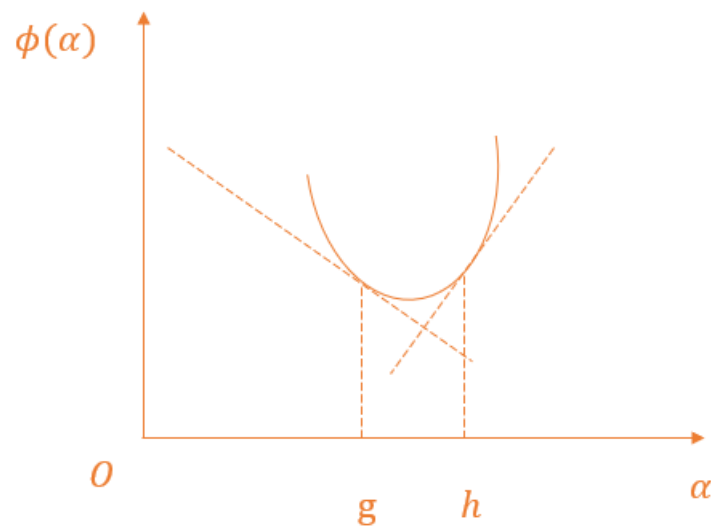
强Wolfe条件使可接受范围的斜率的正值不至于过大，从而将远离极值点的步长排除在外。一般 σ 越小，线搜索越精确，但同时工作量也越大，一般取 $\rho = 0.1, \sigma \in [0.6, 0.8]$



Wolfe-Powell准则的几何意义



强Wolfe-Powell准则的意义



收敛速度问题

➤ 对于一个迭代算法，不仅要求其是收敛的，还要考虑收敛速度问题，一般用阶衡量收敛速度。

线性收敛： 设 \mathbf{x}^* 为 $\min f(\mathbf{x}), \mathbf{x} \in R^n$ 的最优解，由某算法产生的点序列 $\{\mathbf{x}^k\}$ 收敛于 \mathbf{x}^* ，即

$$\lim_{k \rightarrow \infty} \mathbf{x}^k = \mathbf{x}^*$$

若存在一个与 k 无关的常数 $\beta \in (0,1)$ ，对某正数 k_0 ，当 $k > k_0$ 时若有

$$\|\mathbf{x}^{k+1} - \mathbf{x}^*\| \leq \beta \|\mathbf{x}^k - \mathbf{x}^*\|$$

则称序列 $\{\mathbf{x}^k\}$ 为线性收敛，也称该算法线性收敛。

α 阶收敛： 设算法产生序列 $\{\mathbf{x}^k\}$ 收敛于 \mathbf{x}^* ，若存在一个与 k 无关的常数 $\beta \in (0,1)$ 及 $\alpha > 1$ ，对某正数 k_0 ，当 $k > k_0$ 时若有

$$\|\mathbf{x}^{k+1} - \mathbf{x}^*\| \leq \beta \|\mathbf{x}^k - \mathbf{x}^*\|^\alpha$$

则称序列 $\{\mathbf{x}^k\}$ 为 α 阶收敛，也称该算法是 α 阶收敛。

说明： 当 $\alpha = 1$ 时， α 阶收敛就是线性收敛。当 $1 < \alpha < 2$ 时，称算法是超线性收敛，当 $\alpha = 2$ 时，称算法是二阶收敛。阶数越高，收敛速度越快



构造搜索方向问题



➤直接搜索法：在计算过程中只用到目标函数值，不用计算导数

➤解析法：计算过程中要用到目标函数的导数计算



变量轮换法(坐标轴交替下降法)

➤变量轮换法：把多变量函数的优化问题转化为一系列单变量函数的优化问题求解

➤基本思路：搜索方向是各坐标轴的方向，轮流按各坐标轴方向搜索最优点。即从初始点出发，按第 i 个坐标轴方向搜索时，在 n 个变量中，只有 x_i 在变化，其余 $n-1$ 个变量保持不变，依次从 x_1 到 x_n 做 n 次单变量的一维搜索，就完成了变量轮换法的一次迭代

➤设优化问题：

$$\min f(\mathbf{x}), \mathbf{x} \in R^n, f(\mathbf{x}) \in R$$

记 $\mathbf{e}_i = (0, \dots, 1, \dots, 0)^T \quad (i = 1, 2, \dots, n)$

即 \mathbf{e}_i 是第 i 个分量为1其余分量为0的单位分量



变量轮换法

➤ 算法步骤:

Step 1: 给定初始点 $\mathbf{x}^0 = (x_1, x_2, \dots, x_n)^T$, $k=0$;

Step 2: 检查点 \mathbf{x}^k 是否满足终止条件, 满足则结束。否则转Step 3;

Step 3: 完成变量轮换法的一次迭代

记 $\mathbf{y}_0 = \mathbf{x}^k$,
for $i = 1:n$

从 \mathbf{y}_{i-1} 出发, 沿 \mathbf{e}_i 进行线搜索, 记求得的最优步长为 α_i , 可得到新点 \mathbf{y}_i

$$\mathbf{y}_i = \mathbf{y}_{i-1} + \alpha_i \mathbf{e}_i$$

其中 $\alpha_i = \arg \min f(\mathbf{y}_{i-1} + \alpha \mathbf{e}_i)$

Step4: 令 $\mathbf{x}^{k+1} = \mathbf{y}_n$, $k = k + 1$, 转至step 2

优点: 搜索方向的获得不需要花费成本, 基本思想简单、容易理解

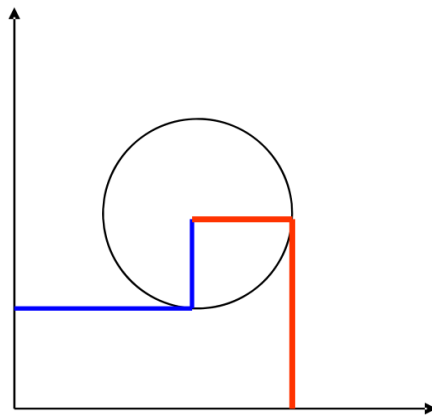
缺点: 收敛速度较慢, 搜索效率较低, 对于一般问题未必收敛。

只适用于具有特殊结构的函数

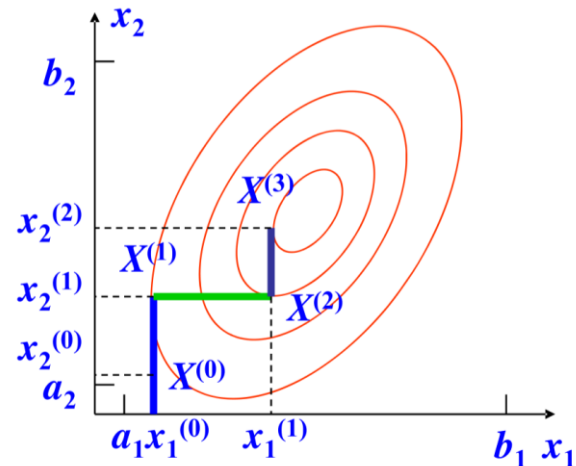


变量轮换法

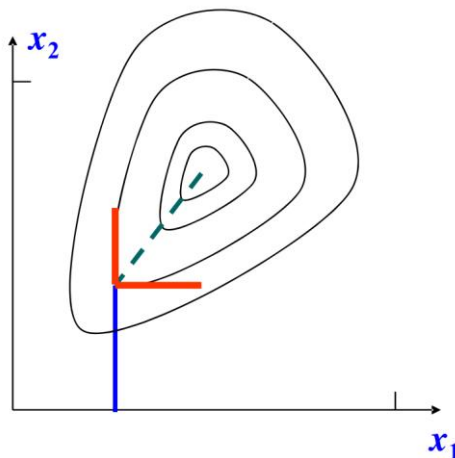
若目标函数的等值线为圆形(二维)、球形(三维)或长短轴平行于坐标轴的椭圆形, 搜索最快。变量之间无交互作用



若目标函数的等值线类似于椭圆形, 且长短轴不平行于坐标轴, 搜索速度较慢。变量之间存在弱交互作用



若目标函数的等值线出现山脊时, 该方法无效。变量之间存在强交互作用



最速下降法(梯度法)

➤考虑无约束优化问题

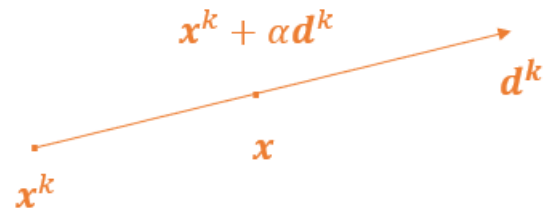
$$\min f(\mathbf{x}), \mathbf{x} \in R^n, f(\mathbf{x}) \in R$$

$f(\mathbf{x})$ 具有一阶连续偏导数, 有极小点 \mathbf{x}^*

设已求得 \mathbf{x}^* 的第 k 次近似值 \mathbf{x}^k , 为了求得 $k+1$ 次近似值 \mathbf{x}^{k+1} , 需选方向 \mathbf{d}^k 。看一下 \mathbf{d}^k 应该具备什么特征?

设 \mathbf{d}^k 已选定, 作射线如图

$$\mathbf{x}^k + \alpha \mathbf{d}^k = \mathbf{x}$$



其中, $\alpha > 0$. $\|\mathbf{d}^k\| = 1$. \mathbf{d}^k 为某个下降方向

现将 $f(\mathbf{x})$ 在 \mathbf{x}^k 点处作一阶泰勒展开:

$$f(\mathbf{x}) = f(\mathbf{x}^k + \alpha \mathbf{d}^k) = f(\mathbf{x}^k) + \alpha \nabla f(\mathbf{x}^k)^T \mathbf{d}^k + O(\|\alpha \mathbf{d}^k\|)$$

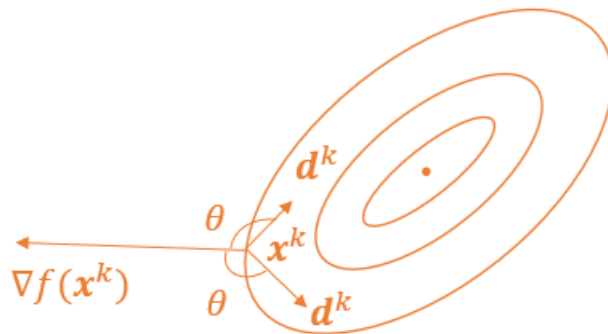
有 $f(\mathbf{x}^k + \alpha \mathbf{d}^k) - f(\mathbf{x}^k) \approx \alpha \nabla f(\mathbf{x}^k)^T \mathbf{d}^k$, 因为 $f(\mathbf{x}^k + \alpha \mathbf{d}^k) - f(\mathbf{x}^k) < 0$

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



最速下降法(梯度法)

- 对于 $\nabla f(\mathbf{x}^k)^T \mathbf{d}^k < 0$ ，即 \mathbf{d}^k 与 \mathbf{x}^k 点处梯度 $\nabla f(\mathbf{x}^k)$ 的点积小于零，或者说两者之间的夹角应大于 90° ，如图所示



- 可以看到，满足这种条件的下降方向有很多，如何选择使目标函数值下降最快的方向？

因为 $\nabla f(\mathbf{x}^k)^T \mathbf{d}^k = \|\nabla f(\mathbf{x}^k)\| \|\mathbf{d}^k\| \cos \theta$ ， θ 为 $\nabla f(\mathbf{x}^k)$ 与 \mathbf{d}^k 之间夹角，显然当 $\theta = 180^\circ$ 时，目标函数值 $f(\mathbf{x}^k)$ 在 \mathbf{x}^k 点附近下降最快，称为负梯度方向，即

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

为最速下降方向

最速下降法(梯度法)

➤确定搜索方向为负梯度方向后，该方向上的所有点可以表示为

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

同时形成确定步长的优化问题

$$f(\mathbf{x}^k + \alpha_k \mathbf{d}^k) = \min_{\alpha} f(\mathbf{x}^k + \alpha \mathbf{d}^k)$$

可以通过所有的线搜索方法求得步长因子 α

基本思想：当前点 \mathbf{x}^k 处选择负梯度方向 $-\nabla f(\mathbf{x}^k)$ 作为搜索方向
最速下降法的步骤：

Step 1: 给定初始点 \mathbf{x}^0 ，终止误差 $\epsilon > 0$ ， $k=0$ ；

Step 2: 判断是否满足终止条件 $\|\nabla f(\mathbf{x}^k)\| < \epsilon$ ，满足则终止；否则转step 3；

Step 3: 构造负梯度方向 $\mathbf{d}^k = -\nabla f(\mathbf{x}^k)$ ；

Step 4: 进行线搜索求得 α_k ，令 $\mathbf{x}^{k+1} = \mathbf{x}^k + \alpha_k \mathbf{d}^k = \mathbf{x}^k - \alpha_k \nabla f(\mathbf{x}^k)$ ， $k = k + 1$ ，转Step 2



最速下降法的优缺点

优点:

1. 思想简单, 每一步迭代只需计算一阶导数, 代价较小
2. 对于well-conditioned, strongly convex问题收敛快

well-conditioned: 函数 $f(\mathbf{x})$ 凸, 一阶可微, 且 $\nabla f(\mathbf{x})$ 是Lipschitz continuous with constant $L > 0$, 即

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 \leq L\|\mathbf{x} - \mathbf{y}\|_2, \text{ for any } \mathbf{x}, \mathbf{y}$$

二阶可微函数, 则 $\nabla^2 f(\mathbf{x}) \preceq L\mathbf{I}$

Strongly convex: $f(\mathbf{x}) - \frac{m}{2}\|\mathbf{x}\|_2^2$ 对于某个 $m > 0$ 是凸函数, 如果函数二阶可微, 则 $\nabla^2 f(\mathbf{x}) \succeq m\mathbf{I}$

缺点:

1. 对于不是well-conditioned, strongly convex问题收敛慢(线性收敛)
2. 不能处理一阶不可微的函数
3. 接近极小点时会出现锯齿(Zigzag)现象

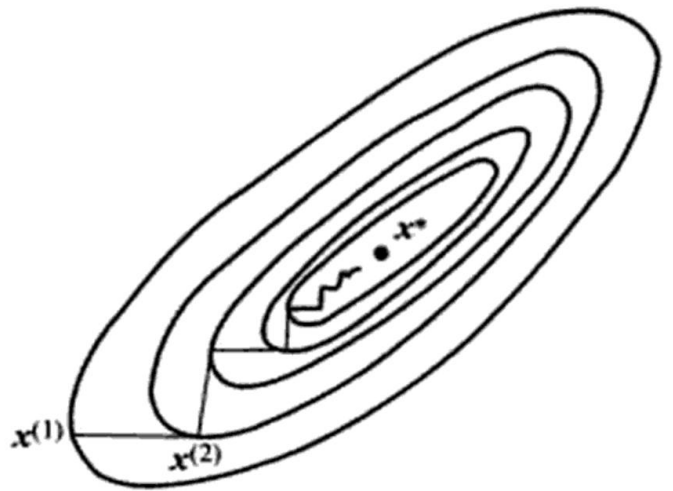


锯齿现象的解释

- 用最速下降法极小化目标函数时，相邻两个搜索方向是正交的
- 若迭代中步长 α_k 是由 $\psi(\alpha) = f(\mathbf{x}^k + \alpha \mathbf{d}^k)$ 的精确最小点，则由 $\psi'(\alpha) = 0$ ，即

$$\begin{aligned}\psi'(\alpha) &= \nabla f(\mathbf{x}^k + \alpha \mathbf{d}^k)^T \mathbf{d}^k \\ &= -\nabla f(\mathbf{x}^{k+1})^T \nabla f(\mathbf{x}^k) \\ &= 0\end{aligned}$$

也就是方向 $\mathbf{d}^{k+1} = -\nabla f(\mathbf{x}^{k+1})$ 与方向 $\mathbf{d}^k = -\nabla f(\mathbf{x}^k)$ 正交。
如此，迭代产生的点列 $\{\mathbf{x}^k\}$ 所走的路径就是“之”字形，如图示



参考文献

- 最优化方法，第4章
- 最优化基础理论与方法，第三章

