

运筹学

(无约束优化)

王焕钢

清华大学自动化系

要点：无约束优化的最优性条件

无约束优化问题 $\min_{X \in R^n} f(X)$

基本假定：目标函数具有二阶导数

梯度 $\nabla f(X) = \frac{\partial f(X)}{\partial X} = \left(\frac{\partial f(X)}{\partial x_1}, \frac{\partial f(X)}{\partial x_2}, \dots, \frac{\partial f(X)}{\partial x_n} \right)^T$

Hesse矩阵

$$\nabla^2 f(X) = \frac{\partial \nabla^T f(X)}{\partial X} = \begin{pmatrix} \frac{\partial^2 f(X)}{\partial x_1 \partial x_1} & \dots & \frac{\partial^2 f(X)}{\partial x_n \partial x_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f(X)}{\partial x_n \partial x_1} & \dots & \frac{\partial^2 f(X)}{\partial x_n \partial x_n} \end{pmatrix}$$

给定方向的二阶泰勒展开

$$f(X + tD) = f(X) + \nabla^T f(X)Dt + \frac{1}{2} D^T \nabla^2 f(X + \xi D)Dt^2$$

最优性条件

1) X^* 是局部最优解的必要条件: $\nabla f(X^*) = 0$

理由: 利用二阶泰勒展开

$$\begin{aligned} f(X^* + tD) - f(X^*) \\ D = -\nabla f(X^*) \quad \Rightarrow \quad &= -\|\nabla^T f(X^*)\|^2 t + \frac{1}{2} D^T \nabla^2 f(X^* + \xi D) D t^2 \\ &= -t \left(\|\nabla^T f(X^*)\|^2 - \frac{1}{2} D^T \nabla^2 f(X^* + \xi D) D t \right) \end{aligned}$$

$$\nabla f(X^*) \neq 0 \quad \Rightarrow \quad f(X^* + tD) - f(X^*) < 0, \quad \forall t \in (0, \varepsilon)$$

$$\Rightarrow X^* \text{ 不是局部最优解}$$

2) X^* 是严格局部最优解的充分条件:

$$\nabla f(X^*) = 0 \quad \nabla^2 f(X^*) > 0$$

理由: $\nabla f(X^*) = 0 \Rightarrow \nabla^T f(X^*)D = 0, \quad \forall D \in R^n$

$$\Rightarrow f(X^* + tD) - f(X^*) = \frac{1}{2} D^T \nabla^2 f(X^* + \xi D) D t^2, \quad \forall D \in R^n$$

$$\nabla^2 f(X^*) > 0 \Rightarrow \nabla^2 f(X^* + \xi D) > 0, \quad \forall \xi \in (0, \hat{\varepsilon})$$

$$\Rightarrow f(X^* + tD) > f(X^*), \quad \forall D \in R^n, t \in (0, \hat{\varepsilon})$$

$$\Rightarrow f(X) > f(X^*), \quad \forall X \in B(X^*, \varepsilon)$$

要点：下降方向法

1847年，法国数学家Cauchy提出**梯度法**

$$\min f(X), X \in \mathbb{R}^n$$

$$\text{迭代算法: } X_{k+1} = X_k + \lambda_k D_k$$

$\lambda_k \in \mathbb{R}^1$ 一维搜索**步长**、 $D_k \in \mathbb{R}^n$ 寻优**方向**

一维**精确**搜索: $X_{k+1} = X_k + \lambda_k^* D_k$

$$f(X_{k+1}) = \min f(X_k + \lambda_k D_k), \lambda_k > 0$$

梯度下降法: $D_k = -\nabla f(X_k)$



奥古斯丁·路易斯·柯西
Augustin Louis Cauchy
1789—1857

基本算法（下降方向法）：

- 1) 任取 $X \in R^n$
- 2) 如果在 X 处找不到下降方向，停止，否则，确定 X 处的下降方向 $D \in R^n$
- 3) 直线搜索确定 t 满足 $f(X + tD) < f(X)$
- 4) 用 $X + tD$ 替换 X ，回到 2) 继续迭代

基本算法（下降方向法）：

- 1) 任取 $X \in R^n$
- 2) 如果在 X 处找不到下降方向，停止，否则，确定 X 处的下降方向 $D \in R^n$
- 3) 直线搜索确定 t 满足 $f(X + tD) < f(X)$
- 4) 用 $X + tD$ 替换 X ，回到 2) 继续迭代

实现算法的**关键**：如何确定下降方向 D ？

要点：梯度下降法

下降方向：负梯度方向

$$D = -\nabla f(X)$$

代入二阶泰勒展开

$$\begin{aligned} f(X + tD) &= f(X) - \|\nabla f(X)\|^2 t + \frac{1}{2} D^T \nabla^2 f(X + \xi D) D t^2 \\ \Rightarrow f(X + tD) - f(X) &= -t \left(\|\nabla f(X)\|^2 - \frac{1}{2} D^T \nabla^2 f(X + \xi D) D t \right) \end{aligned}$$

只要 $\nabla f(X) \neq 0$ ，就有 $\|\nabla f(X)\| > 0$ ，一定存在 $\bar{t} > 0$

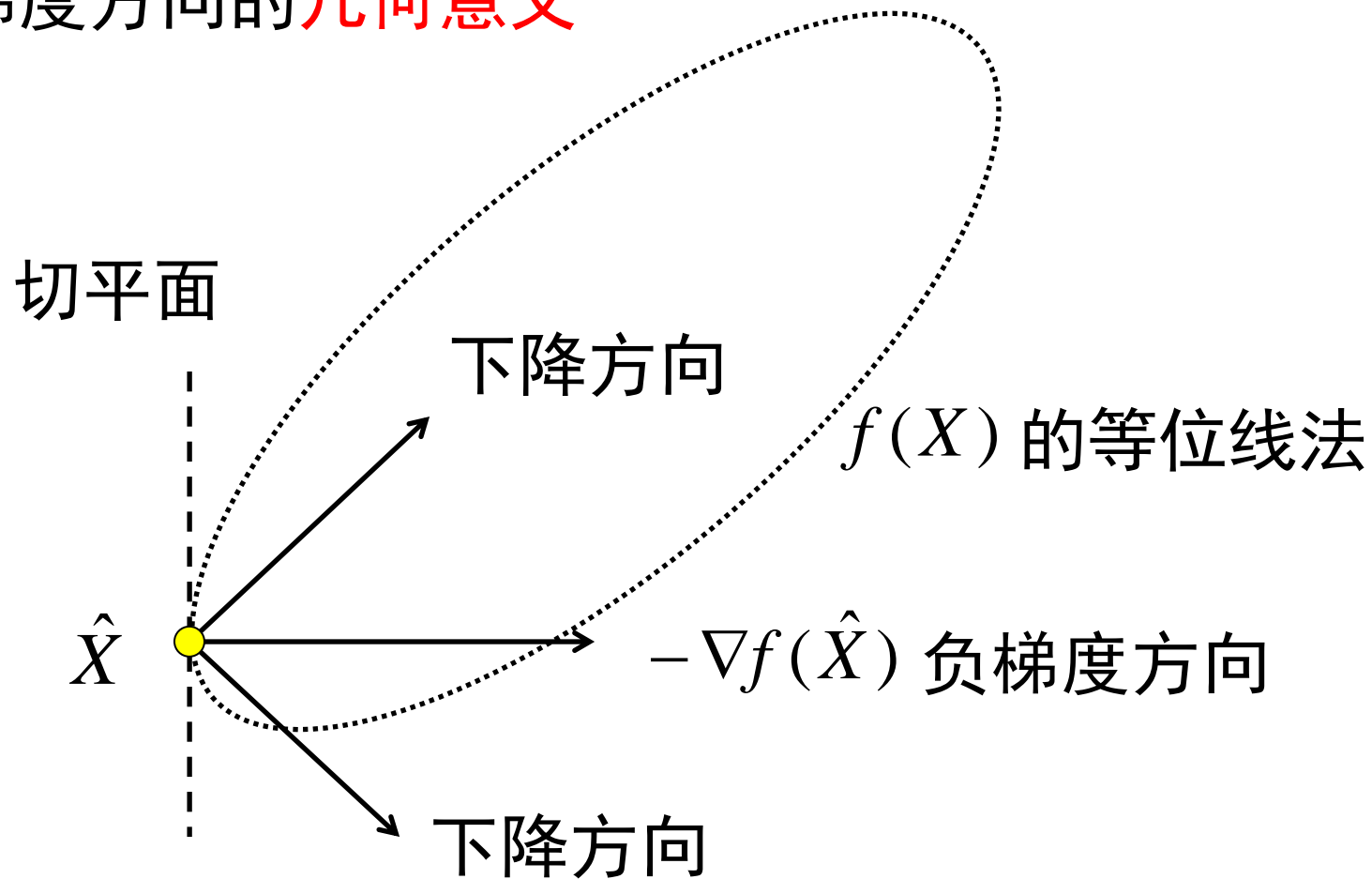
满足 $f(X + tD) < f(X)$, $\forall 0 < t \leq \bar{t}$

所以负梯度方向是下降方向

梯度下降法

- 1) 任取 $\hat{X} \in R^n$
- 2) 计算 $D = -\nabla f(\hat{X})$
- 3) 如果 $\|D\| \leq \delta$ 其中 δ 是预先设定的阈值，停止计算，以 \hat{X} 为所求解，否则进行直线搜索，确定能够满足 $f(\hat{X} + \hat{t}D) < f(\hat{X})$ 的 $\hat{t} > 0$
- 4) 用 $\hat{X} + \hat{t}D$ 替换 \hat{X} ，然后回到 2) 继续迭代

负梯度方向的几何意义



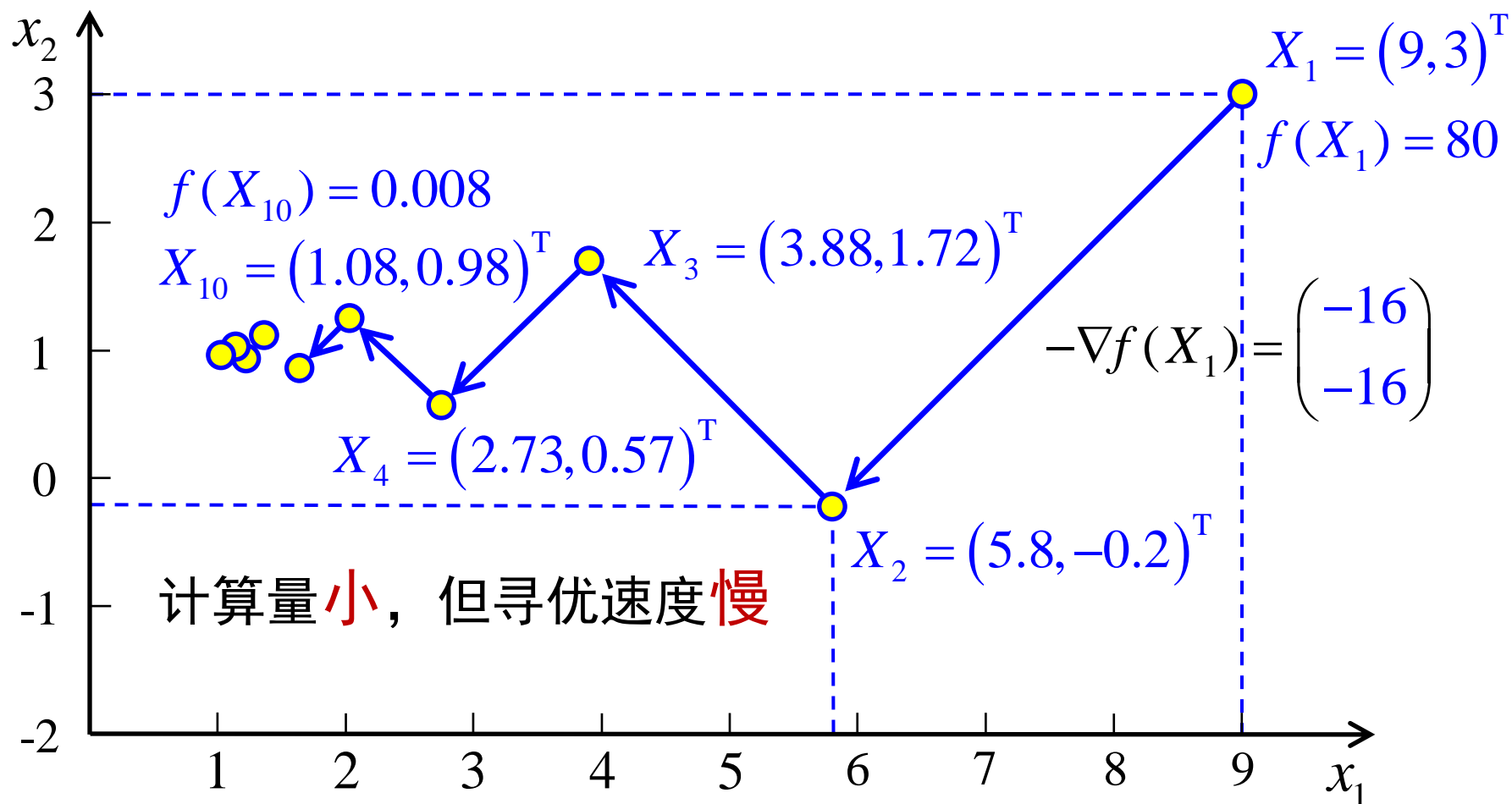
负梯度方向和切平面垂直

梯度下降法的寻优过程

$$\min f(X) = x_1^2 + 4x_2^2 - 2x_1 - 8x_2 + 5$$

梯度下降法

$$X_{k+1} = X_k - \lambda_k^* \nabla f(X_k)$$



要点：负梯度方向的缺陷

一维精确搜索的特性

下降方向 D_1

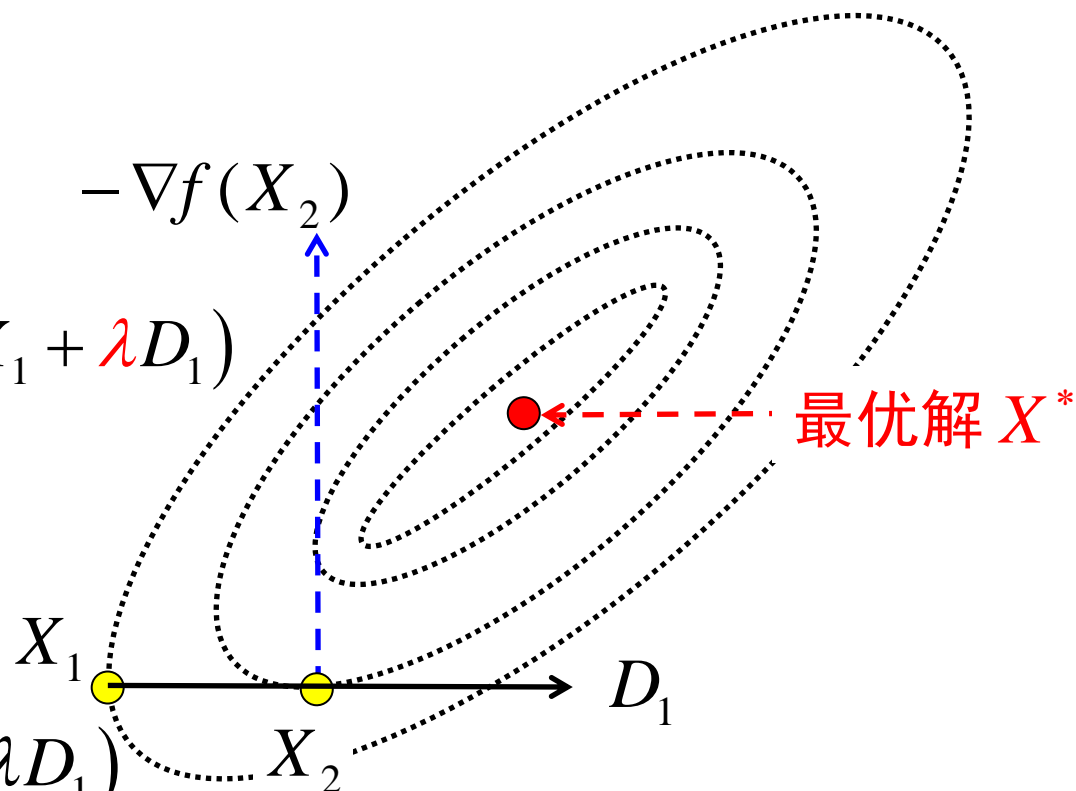
$$X_2 = X_1 + \lambda_1^* D_1$$

λ_1^* 是优化问题 $\min_{\lambda > 0} f(X_1 + \lambda D_1)$ 的最优解

$$\frac{df(X_1 + \lambda D_1)}{d\lambda}$$

$$= \frac{df(X_1 + \lambda D_1)}{dX^T} \frac{d(X_1 + \lambda D_1)}{d\lambda}$$

$$\Rightarrow \nabla^T f(X_1 + \lambda_1^* D_1) D_1 = 0 \quad \Rightarrow \quad \nabla^T f(X_2) D_1 = 0$$



精确搜索得到新点的梯度方向与搜索方向正交

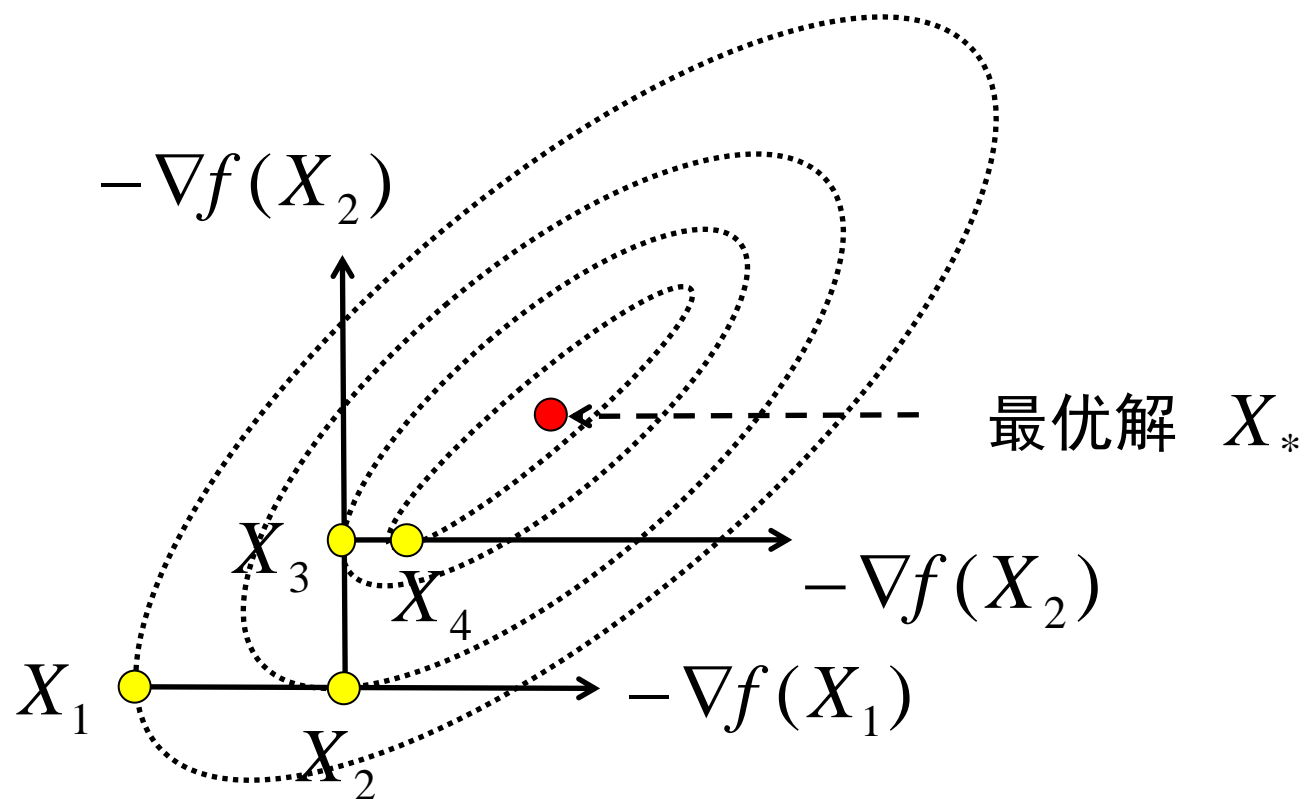
负梯度方向的特点

设 \hat{X}' 是在 \hat{X} 处沿负梯度方向 $D = -\nabla f(\hat{X})$ 进行一维搜索能得到的最好的点，由前面的结果可知

$$\nabla^T f(\hat{X} + \hat{t}D) \nabla f(\hat{X}) = 0$$

即沿负梯度方向精确搜索前进时，相邻两点的梯度互相垂直

负梯度方向的缺陷



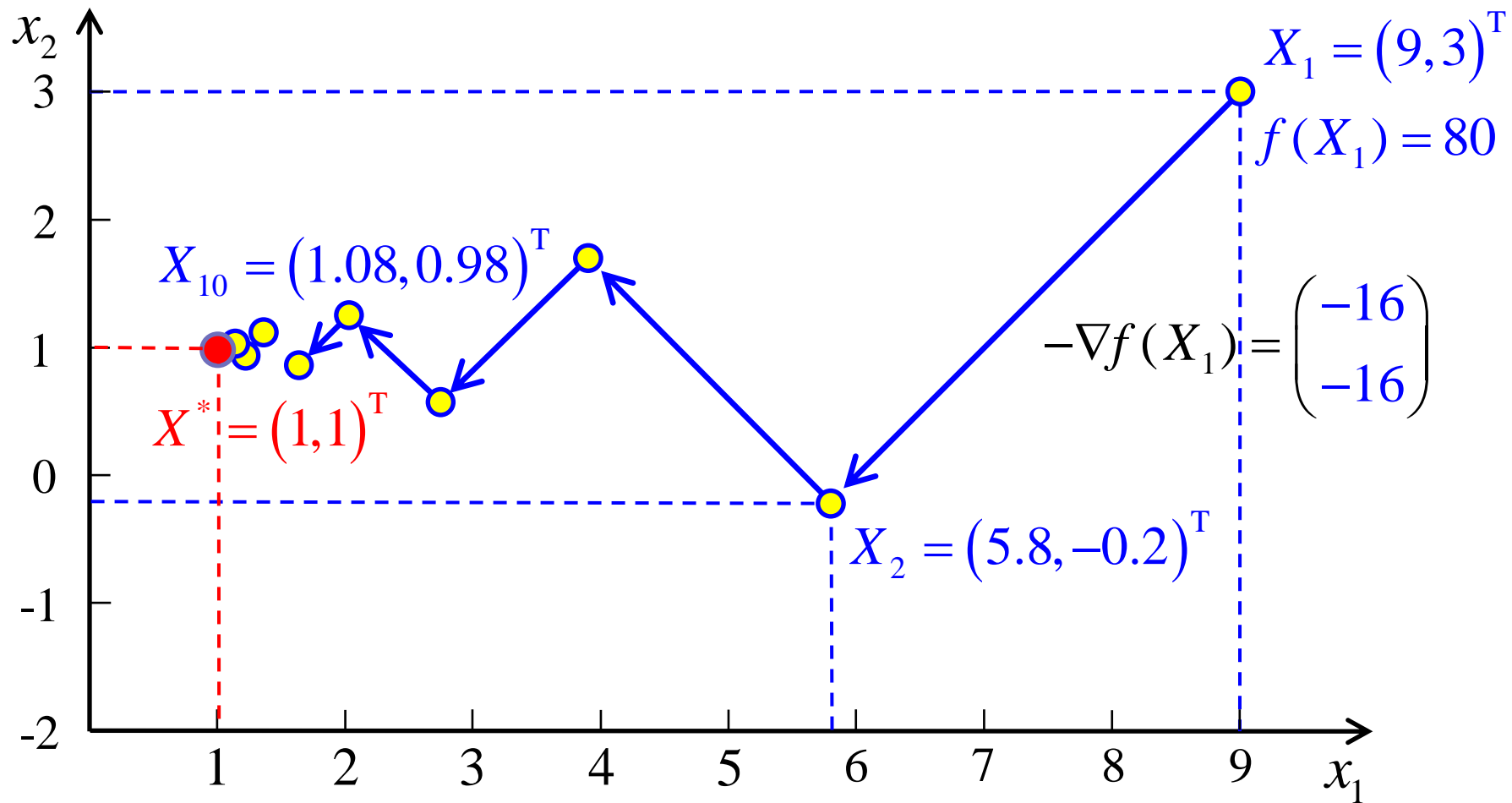
梯度下降法是沿锯齿状路线前进，接近最优解时一维搜索效率很低，前进速度很慢

改进梯度下降法的思路

$$\min f(X) = x_1^2 + 4x_2^2 - 2x_1 - 8x_2 + 5$$

梯度下降法

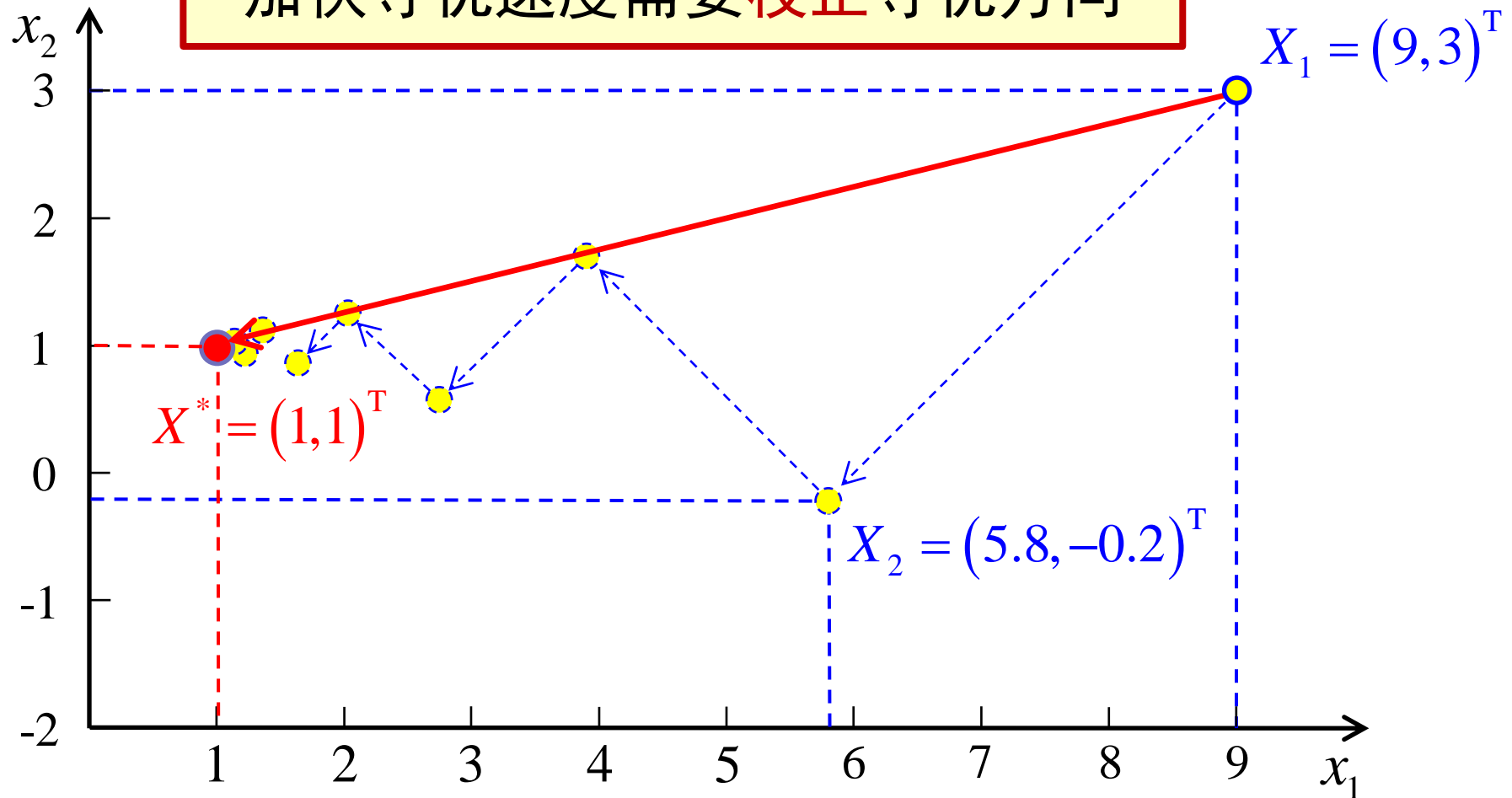
$$X_{k+1} = X_k - \lambda_k^* \nabla f(X_k)$$



改进梯度下降法的思路

$$\min f(X) = x_1^2 + 4x_2^2 - 2x_1 - 8x_2 + 5$$

加快寻优速度需要校正寻优方向



要点：利用梯度方向生成其它下降方向

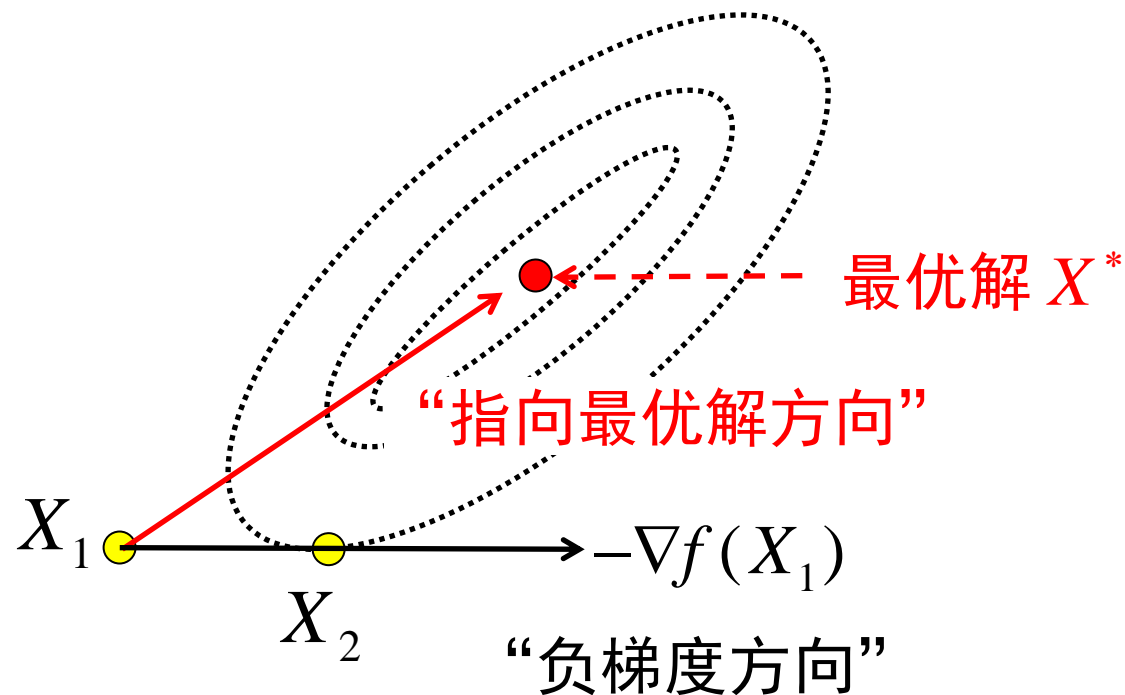
二次函数 $f(X) = \frac{1}{2}(X - X^*)^T A(X - X^*)$ ，其中 $A \in R^{n \times n}$ 为对称正定矩阵， X^* 是固定点

该函数等值面是以 X^* 为中心的椭球面，显然 X^* 为极小值点。

$$X_{k+1} = X_k + \lambda_k D_k$$

沿梯度方向：

$$X_2 = X_1 - \lambda_1 \nabla f(X_1)$$



利用梯度方向生成其它下降方向

任取一个正定矩阵 $Q \in R^{n \times n}$ ，令 $D = -Q\nabla f(X)$

代入二阶泰勒展开可得

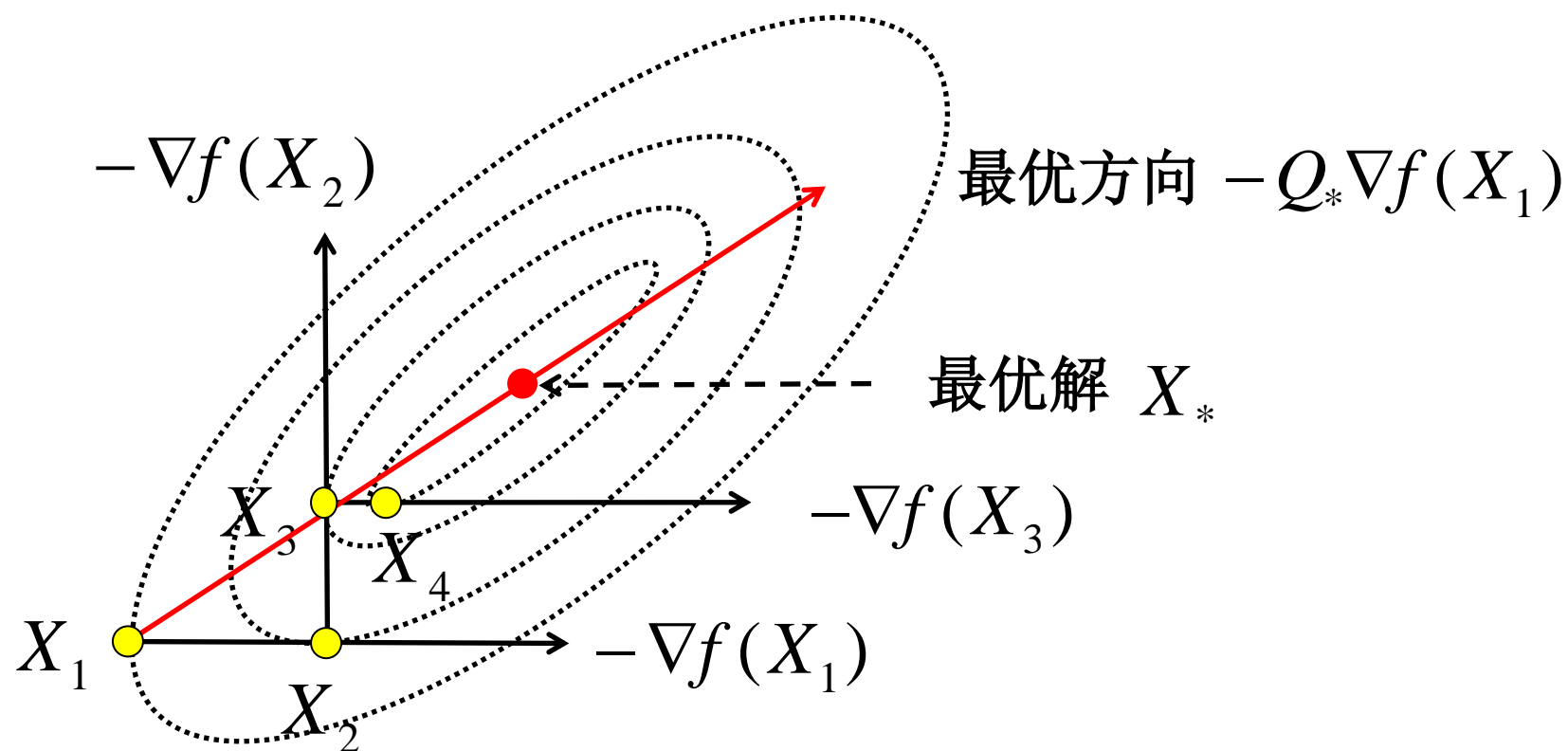
$$\begin{aligned} f(X + tD) - f(X) \\ = -t \left(\nabla^T f(X) Q \nabla f(X) - \frac{1}{2} D^T \nabla^2 f(X + \xi D) D t \right) \end{aligned}$$

只要 $\nabla f(X) \neq 0$ ，就有 $\nabla^T f(X) Q \nabla f(X) > 0$ ，一定存在 $\bar{t} > 0$ 满足

$$f(X + tD) < f(X), \forall 0 < t \leq \bar{t}$$

所以 D 是下降方向

克服负梯度方向缺陷的途径



用适当的正定矩阵（尺度矩阵）乘负梯度方向，其作用是对后者进行**适当的旋转**，以获得更好的方向

要点： 牛顿方向（广义牛顿法）

将二次正定函数 $f(X) = \frac{1}{2}(X - X^*)^T A(X - X^*)$ 改写
为一般形式 $f(X) = \frac{1}{2}X^T AX + B^T X + c$

$$\nabla f(X) = AX + B \Rightarrow X^* = -A^{-1}B$$

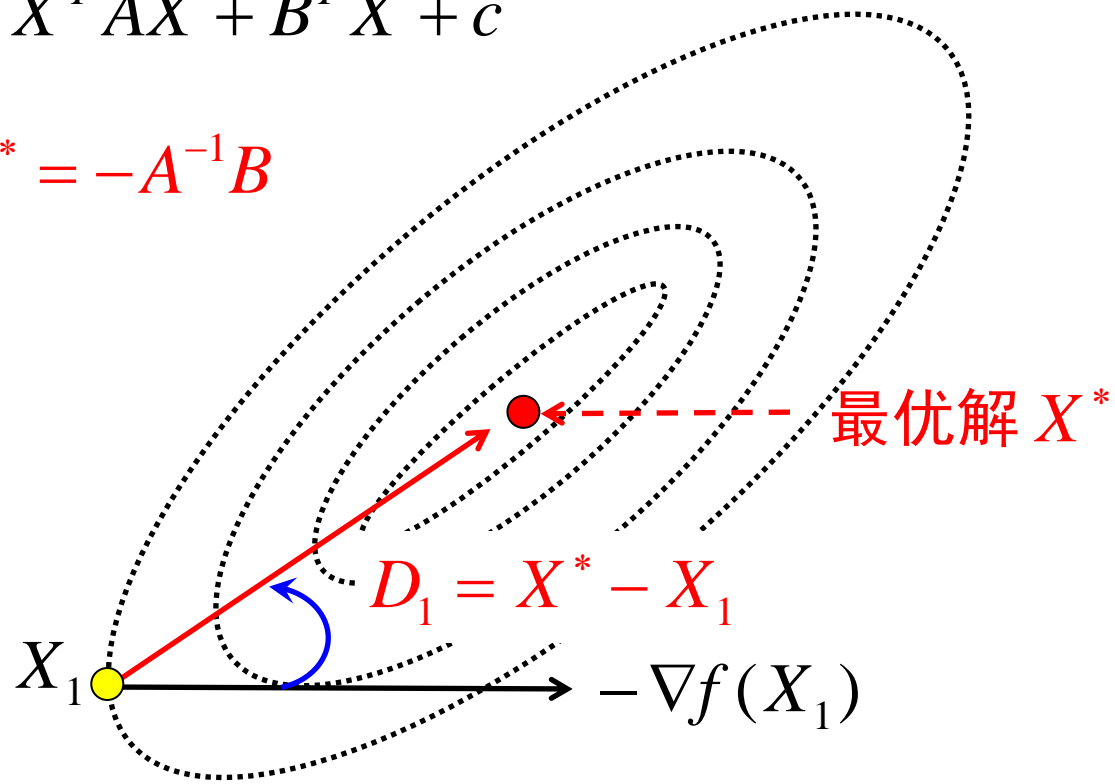
$$\begin{aligned} D_1 &= X^* - X_1 \\ &= -A^{-1}B - X_1 \\ &= -A^{-1}(AX_1 + B) \end{aligned}$$

$$\nabla^2 f(X) = A$$

$$D_1 = -(\nabla^2 f(X_1))^{-1} \nabla f(X_1)$$

$$D_1 = -Q_1 \nabla f(X_1)$$

$$X_k = X_{k-1} - \lambda_{k-1} (\nabla^2 f(X_{k-1}))^{-1} \nabla f(X_{k-1}) \quad \text{牛顿法}$$



正定二次函数的最优方向

对正定二次函数 $f(X) = \frac{1}{2} X^T A X + B^T X + c$

$$\nabla f(X) = AX + B = 0 \Rightarrow X_* = -A^{-1}B$$

$$\nabla^2 f(\hat{X}) = A$$

从任何 \hat{X} 出发, 令

$$D_* = X_* - \hat{X} = -A^{-1}(A\hat{X} + B) = -(\nabla^2 f(\hat{X}))^{-1} \nabla f(\hat{X})$$

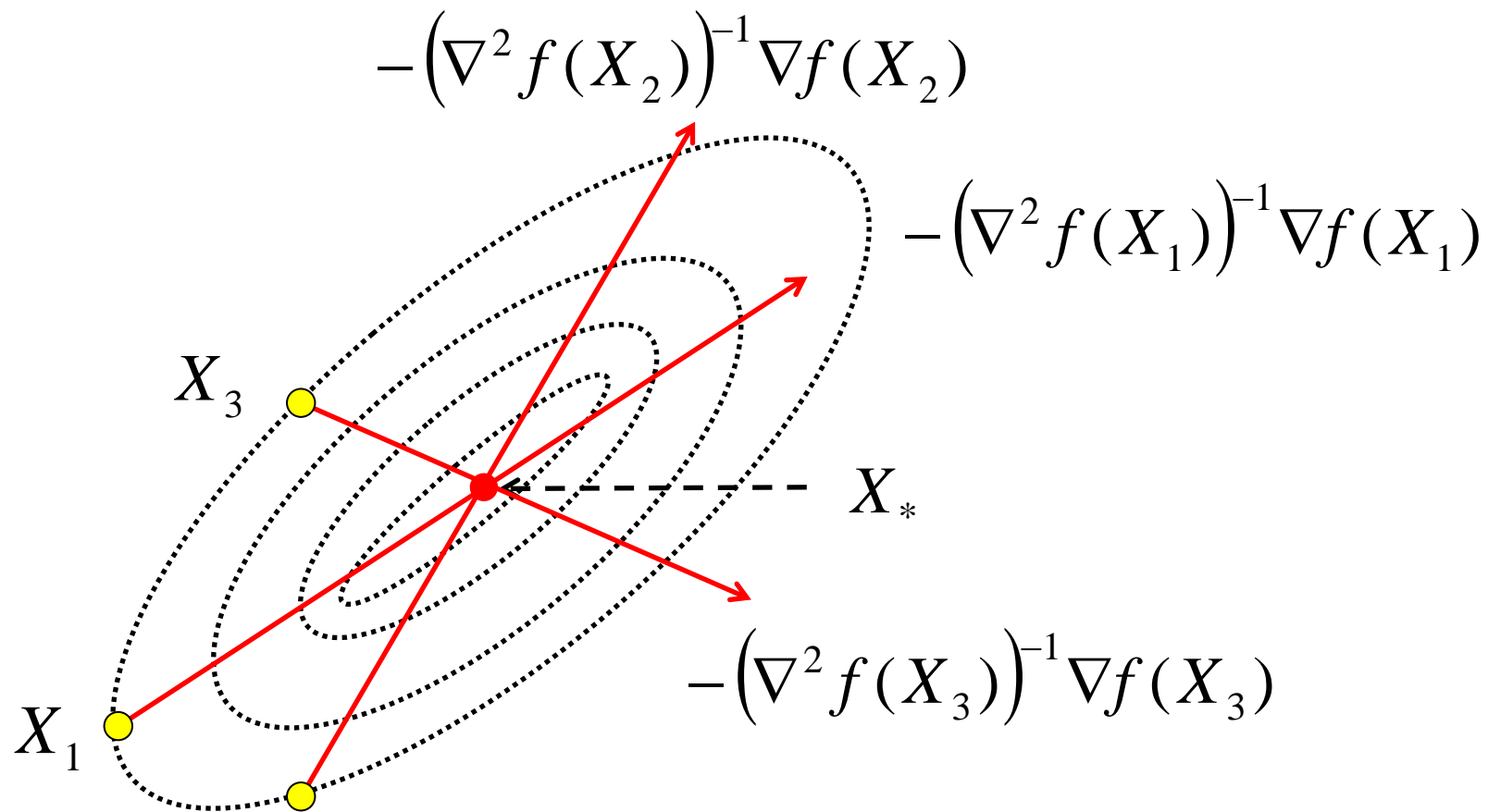
$$t_* = 1$$

显然成立 $\min_{t>0} f(\hat{X} + tD_*) = f(\hat{X} + t_*D_*) = f(X_*)$

说明对任何 \hat{X} , **最优搜索方向**就是

$$D = -(\nabla^2 f(\hat{X}))^{-1} \nabla f(\hat{X}) \quad \text{牛顿方向}$$

正定二次函数的牛顿方向



X_2 说明对任何 \hat{X} ，最优搜索方向就是

$$D = -\left(\nabla^2 f(\hat{X})\right)^{-1} \nabla f(\hat{X})$$

广义牛顿法

- 1) 任取 $\hat{X} \in R^n$
- 2) 如果 $\|\nabla f(\hat{X})\|$ 不大于预先设定的阈值，停止计算，以 \hat{X} 为所求解，否则到下一步
- 3) 计算 $D = -(\nabla^2 f(\hat{X}))^{-1} \nabla f(\hat{X})$ ，进行一维搜索
确定能够满足 $f(\hat{X} + \hat{t}D) < f(\hat{X})$ 的 $\hat{t} > 0$
- 4) 用 $\hat{X} + \hat{t}D$ 替换 \hat{X} ，然后回到 2) 继续迭代

要点：牛顿方向的缺陷

牛顿方向的缺陷

用阻尼牛顿法求解下列问题

$$\min f(x) = x_1^4 + x_1x_2 + (1+x_2)^2$$

初始点 $x^{(1)} = (0,0)^T$ 。

在初始点的梯度和Hessian矩阵分别为

$$\nabla f(x) = \begin{bmatrix} 4x_1^3 + x_2 \\ x_1 + 2(1+x_2) \end{bmatrix} \quad \nabla^2 f(x) = \begin{bmatrix} 12x_1^2 & 1 \\ 1 & 2 \end{bmatrix}$$

$$\nabla f(x^{(1)}) = \begin{bmatrix} 0 \\ 2 \end{bmatrix} \quad \nabla^2 f(x^{(1)}) = \begin{bmatrix} 0 & 1 \\ 1 & 2 \end{bmatrix}$$

在初始点的牛顿方向为

$$d^{(1)} = -\nabla^2 f(x^{(1)})^{-1} \nabla f(x^{(1)}) = \begin{bmatrix} -2 \\ 0 \end{bmatrix}$$

在初始点沿牛顿方向进行一维精确搜索

$$\begin{aligned}\min \varphi(t) &= f\left(x^{(1)} + td^{(1)}\right) \\ &= 16t^4 + 1\end{aligned}$$

可以得到

$$\begin{aligned}\varphi'(t) &= 64t^3 = 0 \\ t^{(1)} &= 0\end{aligned}$$

显然，用阻尼牛顿法不能产生新的点，而初始点并不是无约束优化问题的极小点。

牛顿方向失效的原因在于初始点的Hessian矩阵非正定！

牛顿方向的缺陷

- 1) 每步迭代要计算 $(\nabla^2 f(\hat{X}))^{-1}$ ，计算量大
- 2) $(\nabla^2 f(\hat{X}))^{-1}$ 可能不存在
- 3) $(\nabla^2 f(\hat{X}))^{-1}$ 可能不正定， $D = -(\nabla^2 f(\hat{X}))^{-1} \nabla f(\hat{X})$
不是下降方向

要点：最速下降方向

给定方向的二阶泰勒展开

$$f(X + tD) = f(X) + \nabla^T f(X)Dt + \frac{1}{2}D^T \nabla^2 f(X + \xi D)Dt^2$$

下降方向的充分条件 $\nabla^T f(X)D < 0$

下降方向的必要条件 $\nabla^T f(X)D \leq 0$

最速下降方向

$$\begin{aligned} & \min \left\{ \nabla^T f(X)D \mid \text{s.t. } \|D\| = 1 \right\} \\ \Leftrightarrow & \max \left\{ -\nabla^T f(X)D \mid \text{s.t. } \|D\| = 1 \right\} \Rightarrow \hat{D} \end{aligned}$$

最速下降方向

$$\begin{aligned} & \min \left\{ \nabla^T f(X) D \mid \text{s.t. } \|D\| = 1 \right\} \\ \Leftrightarrow & \max \left\{ -\nabla^T f(X) D \mid \text{s.t. } \|D\| = 1 \right\} \Rightarrow \hat{D} \end{aligned}$$

$$\ell_1 \text{ 范数 } \|D\|_1 = \sum_{i=1}^n |d_i|$$

解决思路：

$$\begin{aligned} -\nabla^T f(X) D &\leq \sum_{i=1}^n \left| \frac{\partial f(X)}{\partial x_i} \right| |d_i| \leq \max_{1 \leq i \leq n} \left| \frac{\partial f(X)}{\partial x_i} \right| \sum_{i=1}^n |d_i| \\ &\Downarrow \\ &\|\nabla f(X)\|_\infty \end{aligned}$$

最速下降方向 $\max \left\{ -\nabla^T f(X)D \mid \text{s.t. } \|D\| = 1 \right\}$

ℓ_1 范数 $\|D\|_1 = \sum_{i=1}^n |d_i|$

$$-\nabla^T f(X)D \leq \|\nabla f(X)\|_\infty \|D\|_1$$

$$\hat{d}_i = \begin{cases} \text{sgn}\left(-\frac{\partial f(X)}{\partial x_i}\right) & \text{if } \left|\frac{\partial f(X)}{\partial x_i}\right| = \|\nabla f(X)\|_\infty \\ 0 & \text{if } \left|\frac{\partial f(X)}{\partial x_i}\right| \neq \|\nabla f(X)\|_\infty \end{cases}$$

$$\nabla f(X)^T \hat{D} = -\|\nabla f(X)\|_\infty$$

$$\ell_p \text{ 范数 } \|D\|_p = \left(\sum_i |d_i|^p \right)^{\frac{1}{p}}, \quad p > 1$$

$$\hat{d}_i = \operatorname{sgn} \left(-\frac{\partial f(X)}{\partial x_i} \right) \left| \frac{\partial f(X)}{\partial x_i} \right|^{q-1} \left(\|\nabla f(X)\|_q \right)^{-\frac{q}{p}}, \quad \forall i$$

$$\nabla f(X)^T \hat{D} = -\|\nabla f(X)\|_q \quad \left(\frac{1}{q} = 1 - \frac{1}{p} \right)$$

$$\ell_\infty \text{ 范数 } \|D\|_\infty = \max_{1 \leq i \leq n} |d_i|$$

$$\hat{d}_i = \operatorname{sgn} \left(-\frac{\partial f(X)}{\partial x_i} \right), \quad \forall i$$

$$\nabla f(X)^T \hat{D} = -\|\nabla f(X)\|_1$$

负梯度方向 (ℓ_2 范数最速下降方向)

$$\hat{D} = -\nabla f(X) \left(\|\nabla f(X)\|_2 \right)^{-1} \Leftrightarrow -\nabla f(X)$$

$$\nabla f(X)^T \hat{D} = -\|\nabla f(X)\|_2$$

牛顿方向 ($\|D\|_{\nabla^2 f(X)} = \left(D^T \nabla^2 f(X) D \right)^{\frac{1}{2}}$ 的最速下降方向)

$$\hat{D} = \frac{-\left(\nabla^2 f(X) \right)^{-1} \nabla f(X)}{\left(\nabla f(X)^T \left(\nabla^2 f(X) \right)^{-1} \nabla f(X) \right)^{\frac{1}{2}}} \Leftrightarrow -\left(\nabla^2 f(X) \right)^{-1} \nabla f(X)$$

$$\nabla f(X)^T \hat{D} = -\nabla f(X)^T \left(\nabla^2 f(X) \right)^{-1} \nabla f(X)$$

要点：共轭梯度方向

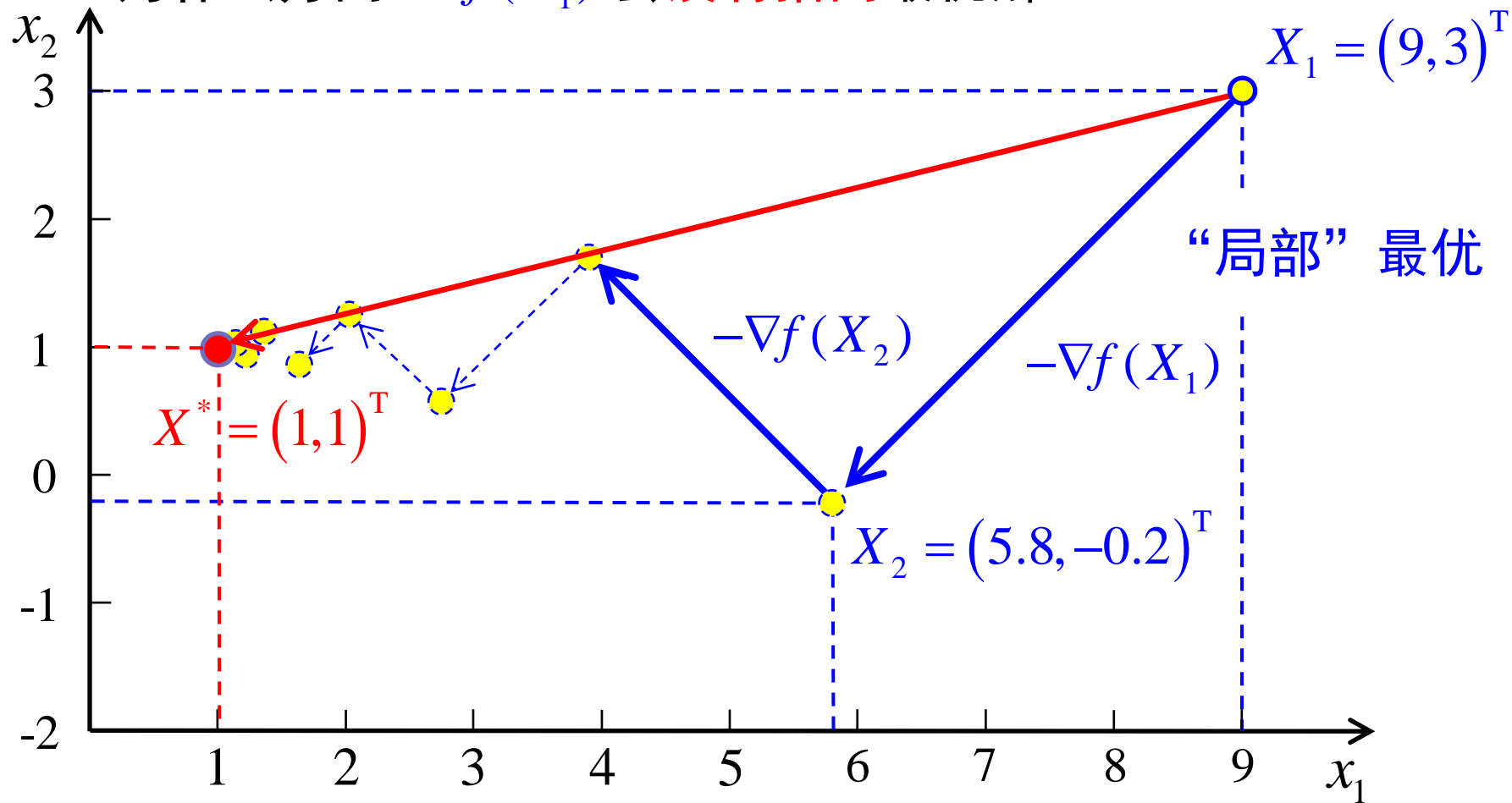
改进梯度下降法的思路

柯西：梯度下降法

$$\min f(X) = x_1^2 + 4x_2^2 - 2x_1 - 8x_2 + 5$$

$$X_{k+1} = X_k - \lambda_k^* \nabla f(X_k)$$

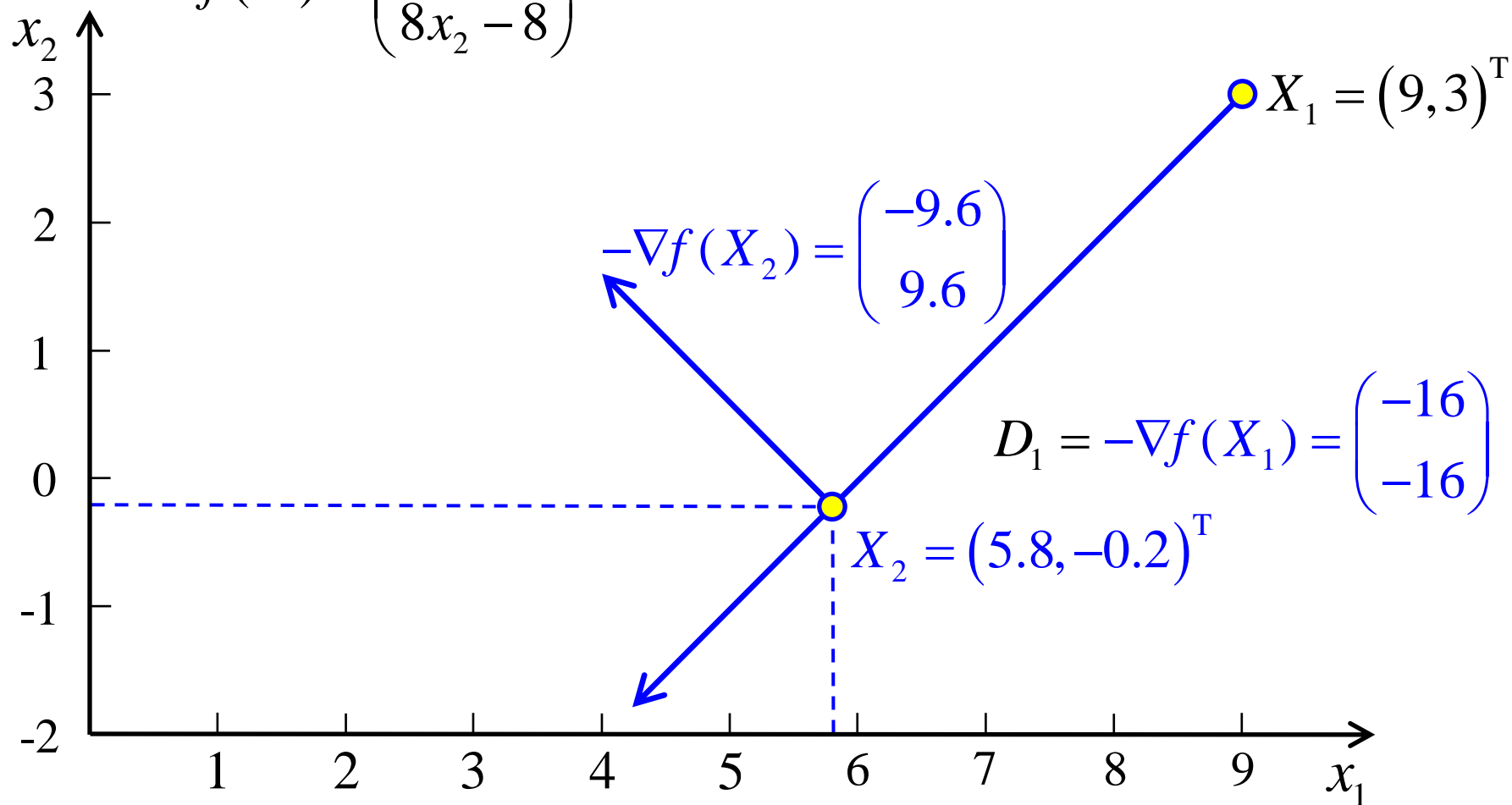
为什么方向 $-\nabla f(X_1)$ 会没有指向最优解 X^* ？



F-R 共轭梯度法 —— 利用共轭梯度寻优

$\min f(X) = x_1^2 + 4x_2^2 - 2x_1 - 8x_2 + 5$ 第一步沿负梯度寻优

$$\nabla f(X) = \begin{pmatrix} 2x_1 - 2 \\ 8x_2 - 8 \end{pmatrix}$$

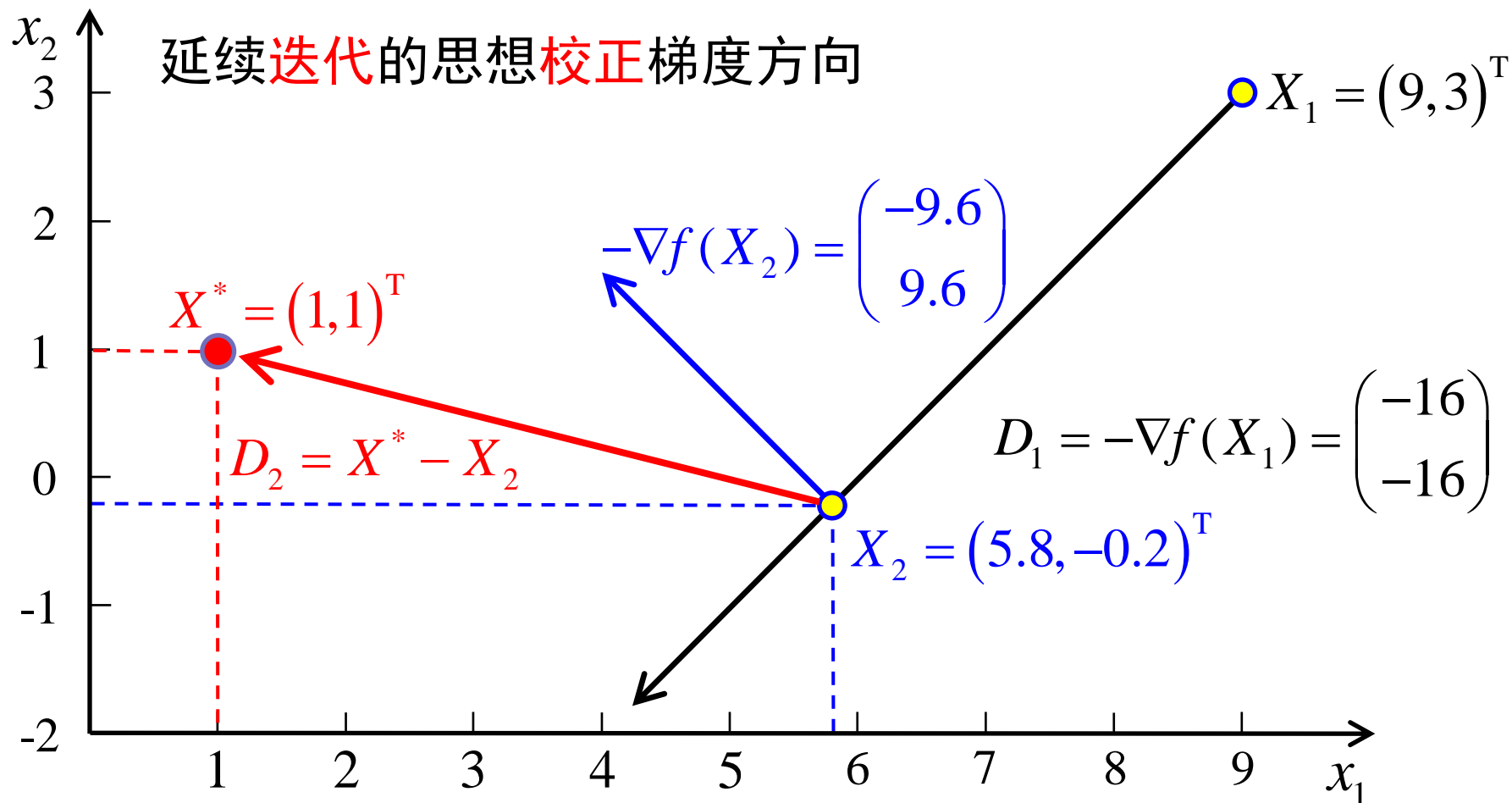


F-R 共轭梯度法 —— 利用共轭梯度寻优

$$\min f(X) = x_1^2 + 4x_2^2 - 2x_1 - 8x_2 + 5$$

$$D_2 = -\nabla f(x_2) + \alpha_1 D_1$$

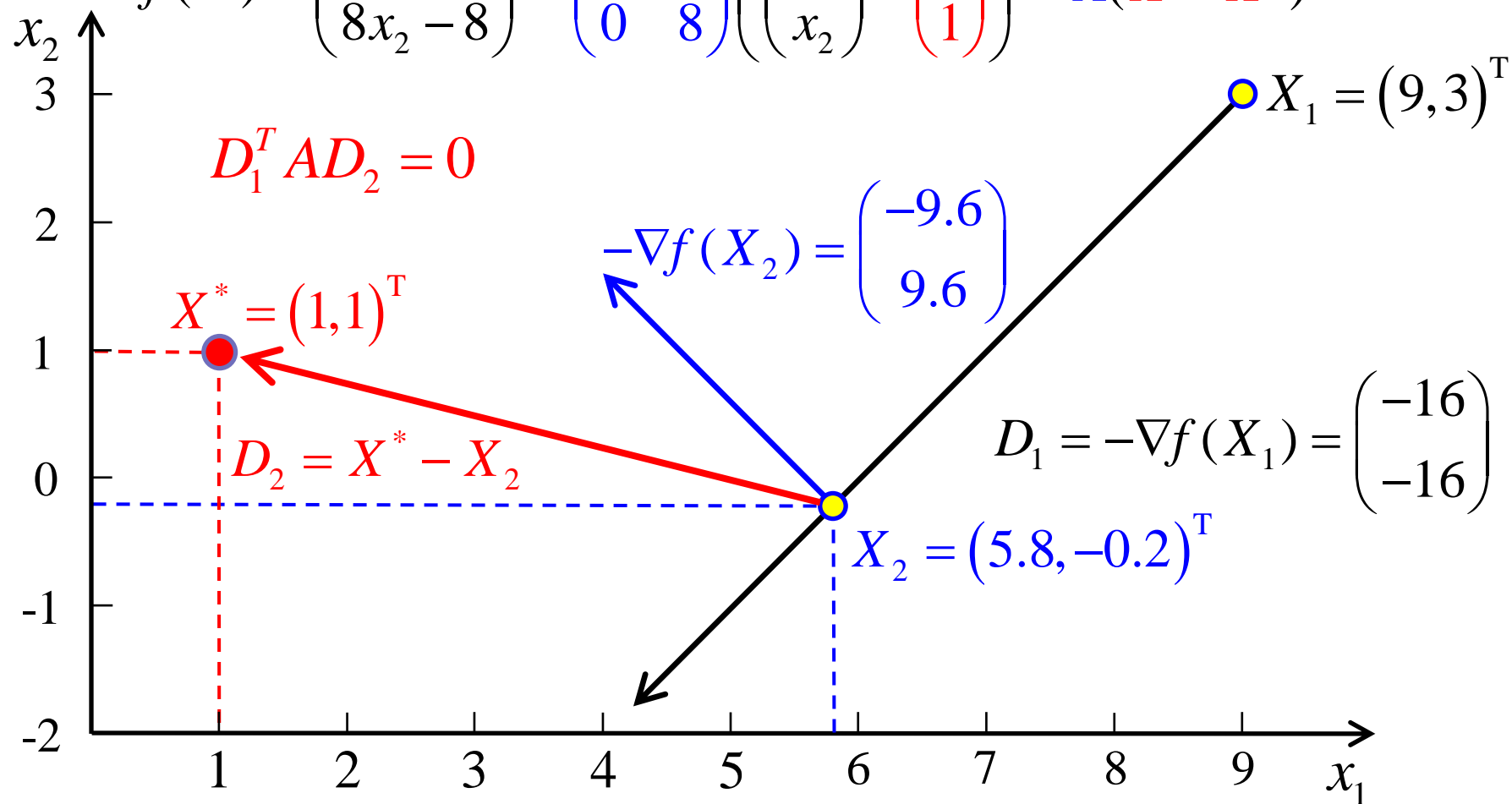
$$X_{k+1} = X_k + \lambda_k D_k$$



F-R 共轭梯度法 —— 利用共轭梯度寻优

$$\min f(X) = x_1^2 + 4x_2^2 - 2x_1 - 8x_2 + 5$$

$$\nabla f(X) = \begin{pmatrix} 2x_1 - 2 \\ 8x_2 - 8 \end{pmatrix} = \begin{pmatrix} 2 & 0 \\ 0 & 8 \end{pmatrix} \left(\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} - \begin{pmatrix} 1 \\ 1 \end{pmatrix} \right) = A(X - X^*)$$



共轭方向法原理之一

共轭方向定义： $A \in R^{n \times n}$ **对称矩阵**， $\vec{p}, \vec{q} \in R^n$ **非零向量**，若 $\vec{p}^T A \vec{q} = 0$ ，称 \vec{p}, \vec{q} 为 A 共轭方向

共轭方向线性无关性

若 $\vec{p}_0, \vec{p}_1, \dots, \vec{p}_{n-1}$ 互为 $A > 0$ 的共轭方向，则它们线性无关

$$\text{理由： } \alpha_0 \vec{p}_0 + \alpha_1 \vec{p}_1 + \dots + \alpha_{n-1} \vec{p}_{n-1} = 0$$

$$\Rightarrow \alpha_0 \vec{p}_0^T A \vec{p}_0 + \alpha_1 \vec{p}_0^T A \vec{p}_1 + \dots + \alpha_{n-1} \vec{p}_0^T A \vec{p}_{n-1} = 0$$

$$\Rightarrow \alpha_k \vec{p}_k^T A \vec{p}_k = 0$$

$$\Rightarrow \alpha_k = 0$$

共轭梯度方向 $D_k = -\nabla f(X_k) + \alpha_{k-1} D_{k-1}$

对于二次正定函数 $f(X) = \frac{1}{2}(X - X^*)^T A(X - X^*)$

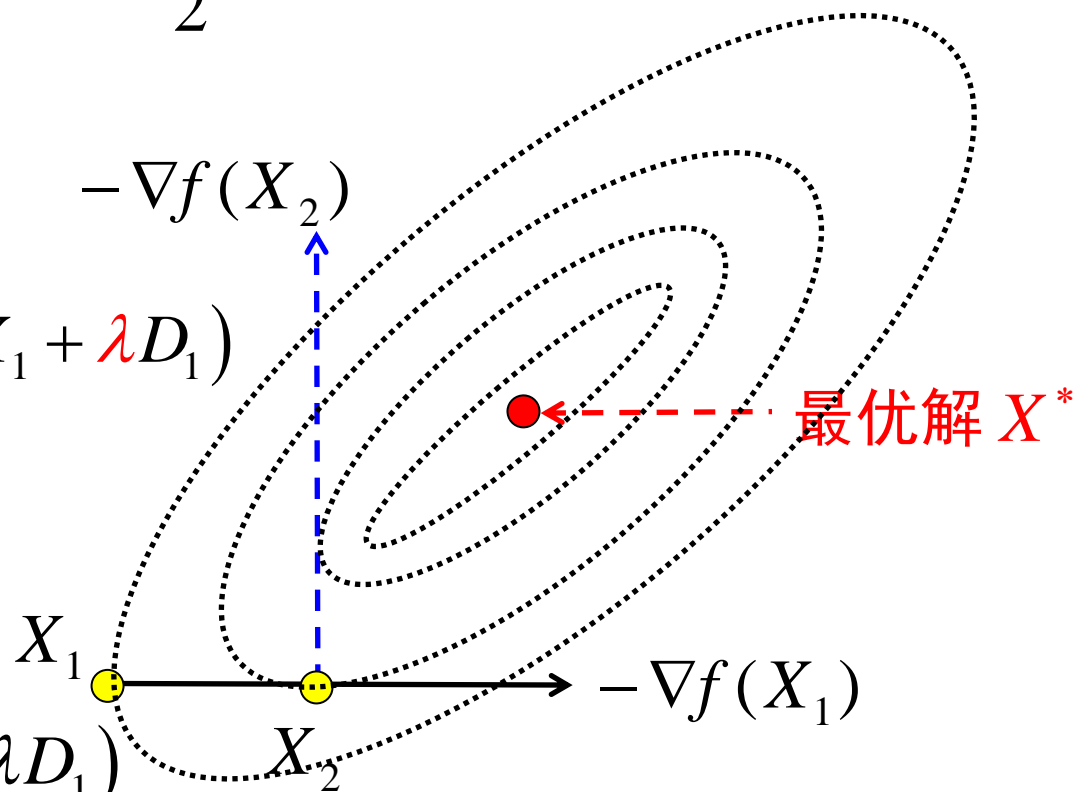
$$D_1 = -\nabla f(X_1)$$

$$X_2 = X_1 + \lambda_1^* D_1$$

λ_1^* 是优化问题 $\min_{\lambda > 0} f(X_1 + \lambda D_1)$
的最优解

$$\begin{aligned} & \frac{df(X_1 + \lambda D_1)}{d\lambda} \\ &= \frac{df(X_1 + \lambda D_1)}{dX^T} \frac{d(X_1 + \lambda D_1)}{d\lambda} \end{aligned}$$

$$\Rightarrow \nabla^T f(X_1 + \lambda_1^* D_1) D_1 = 0 \quad \Rightarrow \quad \nabla^T f(X_2) D_1 = 0$$



$$f(X) = \frac{1}{2}(X - X^*)^T A(X - X^*) \quad D_k = -\nabla f(X_k) + \alpha_{k-1} D_{k-1}$$

$$D_1 = -\nabla f(X_1)$$

$$D_2 = -\nabla f(X_2) + \alpha_1 D_1$$

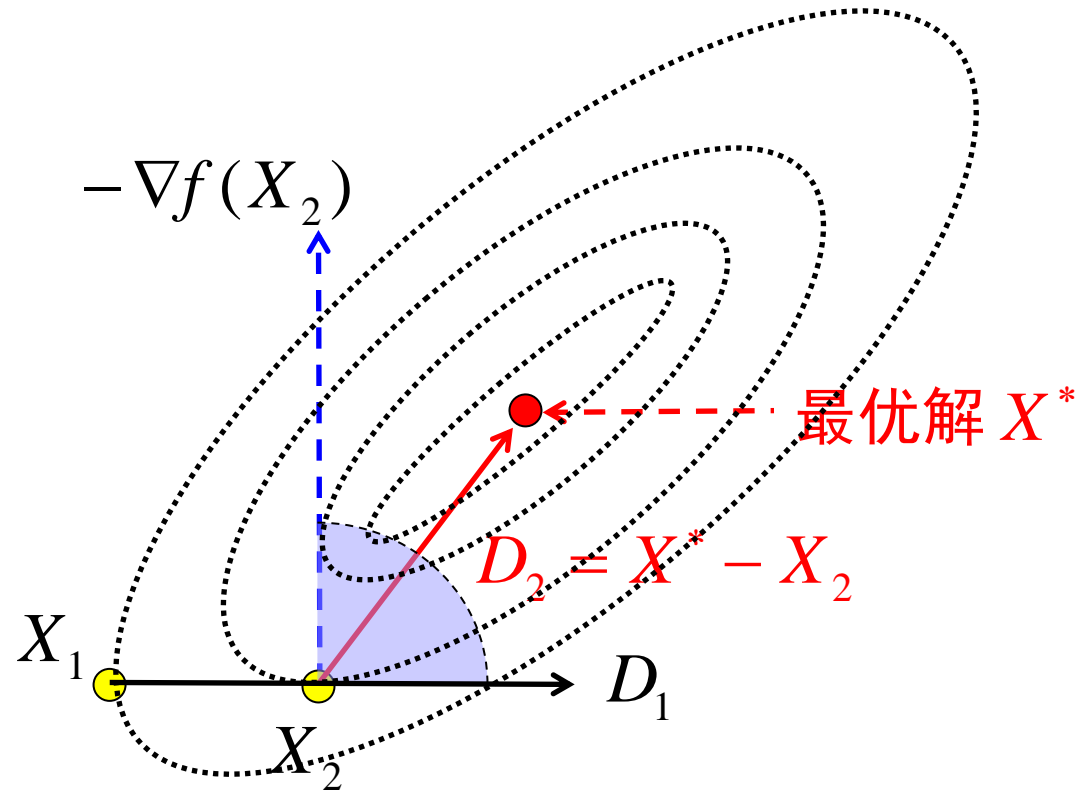
$$\begin{aligned} \nabla f(X_2) &= A(X_2 - X^*) \\ &= -A D_2 \end{aligned}$$

$$D_1^T \nabla f(X_2) = 0$$

$$\Rightarrow -D_1^T A D_2 = 0$$

$$\Rightarrow D_1^T A D_2 = 0$$

D_1 与 D_2 为 A 的共轭方向！



要点：F-R共轭梯度法

F-R 共轭梯度法 —— 利用共轭梯度寻优

$$\min f(X) = x_1^2 + 4x_2^2 - 2x_1 - 8x_2 + 5$$

$$A = \begin{pmatrix} 2 & 0 \\ 0 & 8 \end{pmatrix}$$

“长记性” 寻优方向 D_2 和 D_1 是两个 A 的共轭方向

$$D_2 = -\nabla f(x_2) + \alpha_1 D_1$$

共轭梯度

只需要解决如何计算出合适的参数 α_1

1952年Hestenes和Stiefel提出利用共轭梯度

求解线性方程组 $AX = b, X \in \mathbb{R}^n$

$$\min (X^T AX - b^T X), X \in \mathbb{R}^n$$

Fletcher: 用“简单”解决“复杂”



R. Fletcher 英国
皇家科学院院士

F-R 共轭梯度法 —— F-R共轭梯度法

1964年, Fletcher 和 Reeves提出了适用于一般无约束最优化问题的求解方法: **F-R 共轭梯度法**

梯度下降法

$$X_{k+1} = X_k + \lambda_k^* D_k$$

$$X_{k+1} = X_k + \lambda_k^* D_k$$

$$D_k = \begin{cases} -\nabla f(X_k) & k = 1 \\ -\nabla f(X_k) + \alpha_{k-1} D_{k-1} & k \geq 2 \end{cases}$$

$$D_k = -\nabla f(X_k)$$

相邻两步寻优方向**共轭性** $D_k^T A D_{k-1} = 0$ 和**精确搜索**的特点

$$\alpha_k = \frac{\|\nabla f(X_{k+1})\|^2}{\|\nabla f(X_k)\|^2}$$

F-R 共轭梯度法计算简单、寻优速度快, 在国际上开启了**共轭梯度法**求解非线性规划的研究先河!

要点：参数 α 的计算

$$f(X) = \frac{1}{2}(X - X^*)^T A(X - X^*) \quad D_k = -\nabla f(X_k) + \alpha_{k-1} D_{k-1}$$

$$D_1 = -\nabla f(X_1)$$

$$D_2 = -\nabla f(X_2) + \alpha_1 D_1$$

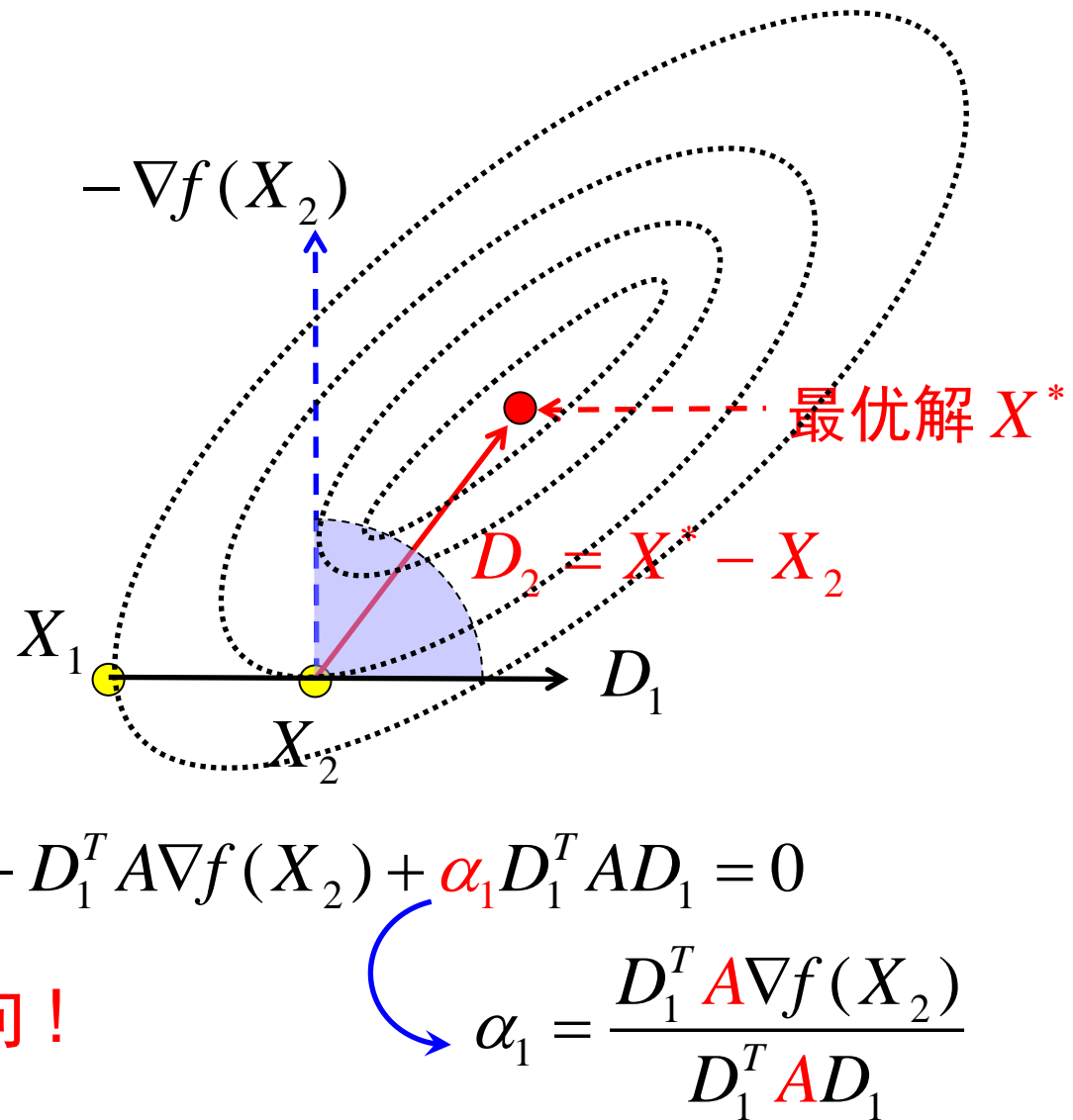
$$\begin{aligned} \nabla f(X_2) &= A(X_2 - X^*) \\ &= -AD_2 \end{aligned}$$

$$D_1^T \nabla f(X_2) = 0$$

$$\Rightarrow -D_1^T AD_2 = 0$$

$$\Rightarrow D_1^T AD_2 = 0 \Rightarrow -D_1^T A \nabla f(X_2) + \alpha_1 D_1^T AD_1 = 0$$

D_1 与 D_2 为 A 的共轭方向！



参数 α 中矩阵 A 的消除方法

由 $X_2 = X_1 + \lambda_1^* D_1 \Rightarrow D_1 = (X_2 - X_1) / \lambda_1^*$

$$\begin{aligned}\alpha_1 &= \frac{D_1^T A \nabla f(X_2)}{D_1^T A D_1} = \frac{\nabla^T f(X_2) A D_1}{D_1^T A D_1} \\&= \frac{\nabla^T f(X_2) A (X_2 - X_1) / \lambda_1^*}{D_1^T A (X_2 - X_1) / \lambda_1^*} = \frac{\nabla^T f(X_2) A (X_2 - X_1)}{D_1^T A (X_2 - X_1)} \\&= \frac{\nabla^T f(X_2) (\nabla f(X_2) - \nabla f(X_1))}{D_1^T (\nabla f(X_2) - \nabla f(X_1))} = \frac{\nabla^T f(X_2) (\nabla f(X_2) + D_1)}{D_1^T (\nabla f(X_2) + D_1)} \\&= \frac{\nabla^T f(X_2) \nabla f(X_2)}{D_1^T D_1} = \frac{\|\nabla f(X_2)\|^2}{\|\nabla f(X_1)\|^2}\end{aligned}$$

要点：F-R共轭梯度法计算示例

F-R 共轭梯度法 —— 寻优速度对比

$$X_{k+1} = X_k + \lambda_k^* D_k$$

$$\min f(X) = x_1^2 + 4x_2^2 - 2x_1 - 8x_2 + 5$$

$$D_k = \begin{cases} -\nabla f(X_k) & k=1 \\ -\nabla f(X_k) + \alpha_{k-1} D_{k-1} & k \geq 2 \end{cases}$$

$$\alpha_k = \frac{\|\nabla f(X_{k+1})\|^2}{\|\nabla f(X_k)\|^2}$$

$$X_1 = (9, 3)^T$$

F-R 法计算步骤

① $D_1 = -\nabla f(X_1)$

② $\min_{\lambda_1 > 0} f(X_1 + \lambda_1 D_1), X_2 = X_1 + \lambda_1^* D_1$

③ $\alpha_1 = \frac{\|\nabla f(X_2)\|^2}{\|\nabla f(X_1)\|^2}$

④ $D_2 = -\nabla f(X_2) + \alpha_1 D_1$

⑤ $\min_{\lambda_2 > 0} f(X_2 + \lambda_2 D_2), X_3 = X_2 + \lambda_2^* D_2$

计算结果

$$D_1 = -(16, 16)^T$$

$$\lambda_1^* = 0.2, X_2 = (5.8, -0.2)^T$$

$$\nabla f(X_2) = (9.6, -9.6)^T$$

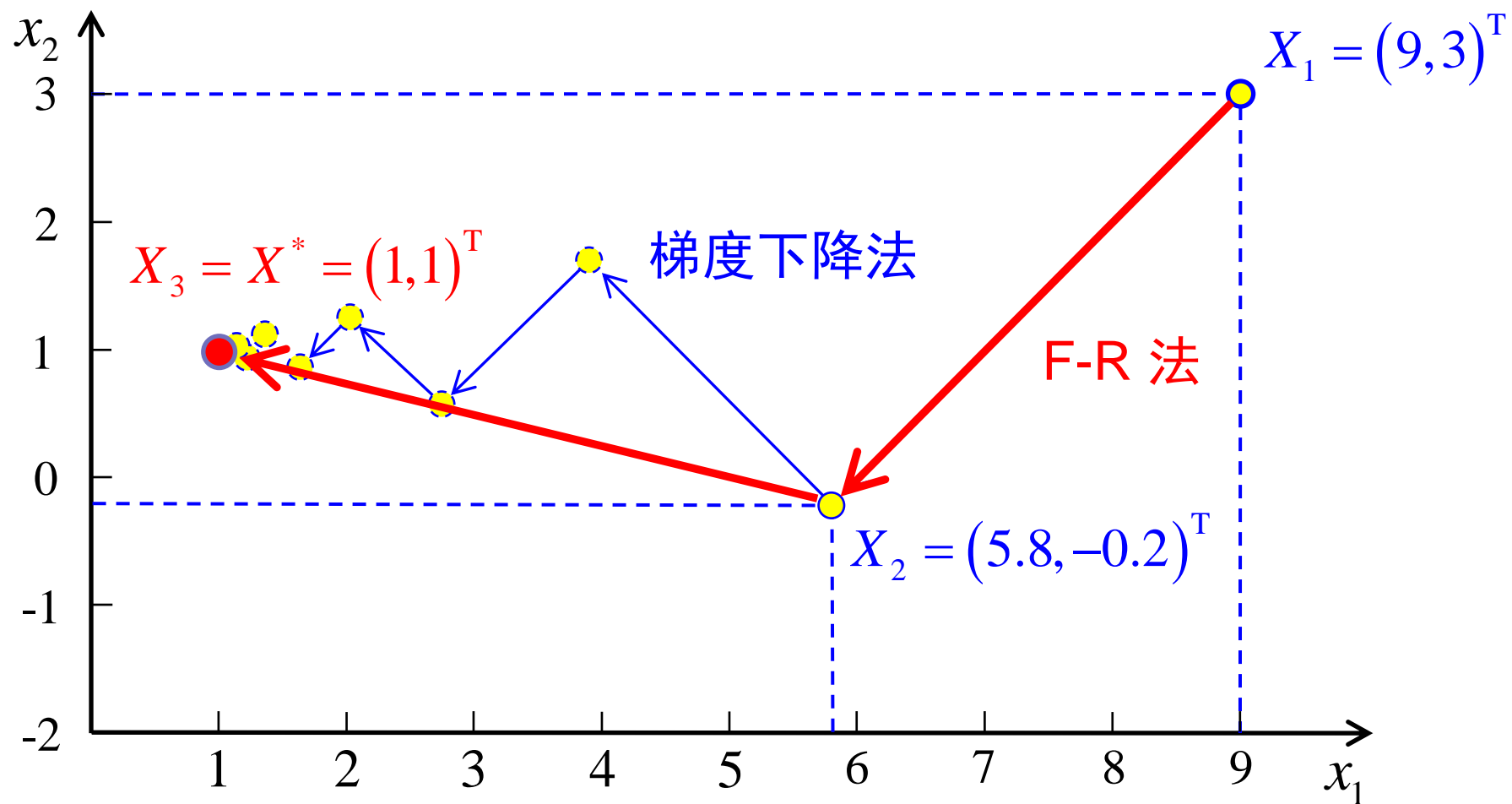
$$\alpha_1 = 0.36$$

$$D_2 = (-15.36, 3.84)^T$$

$$\lambda_2^* = 0.3125, X_3 = (1, 1)^T$$

F-R 共轭梯度法 —— 寻优轨迹对比

$$\min f(X) = x_1^2 + 4x_2^2 - 2x_1 - 8x_2 + 5$$



要点：与一维最优解的梯度的正交性

共轭方向和一维最优解的梯度的正交性

条件: $f(X) = 0.5X^TAX + B^TX + C, \quad A > 0$

$\vec{p}_0, \vec{p}_1, \dots, \vec{p}_{n-1}$ 为 A 的共轭方向

$X_0 \in R^n$ 是任意的出发点

由下述一维搜索依次确定 X_1, X_2, \dots, X_n

$$f(X_{k+1}) = f(X_k + t_k \vec{p}_k) = \min_{t \in R} f(X_k + t \vec{p}_k)$$
$$k = 0, 1, \dots, n-1$$

结论: $\vec{p}_j^T \nabla f(X_k) = 0, \quad \forall 0 \leq j < k$

理由: $\min_{t>0} f(X_k + t\vec{p}_k) \Rightarrow t_k = -\frac{\vec{p}_k^T \nabla f(X_k)}{\vec{p}_k^T A \vec{p}_k}, \forall 0 \leq k \leq n-1$

$$X_k = X_{k-1} + t_{k-1} \vec{p}_{k-1} \Rightarrow X_k = X_0 + \sum_{i=0}^{k-1} t_i \vec{p}_i$$

$$\Rightarrow \nabla f(X_k) = \nabla f(X_0) + \sum_{i=0}^{k-1} t_i A \vec{p}_i$$

$$\begin{aligned} \vec{p}_j^T \nabla f(X_k) &= \vec{p}_j^T \nabla f(X_0) + t_j \vec{p}_j^T A \vec{p}_j \\ &= \vec{p}_j^T \nabla f(X_0) - \vec{p}_j^T \nabla f(X_j), \quad \forall 0 \leq j < k \\ \Rightarrow \end{aligned}$$

$$\vec{p}_j^T \nabla f(X_j) = \vec{p}_j^T \nabla f(X_0), \quad \forall j$$

$$\Rightarrow \vec{p}_j^T \nabla f(X_k) = 0, \quad \forall 0 \leq j < k$$

推论：沿共轭方向寻优的每个 $X_k, k = 1, 2, \dots, n$ 都满足

$$f(X_k) = \min \left\{ f(X) \mid \text{s.t. } X = X_0 + \sum_{j=0}^{k-1} \beta_j \vec{p}_j \right\}$$

理由： $X = X_0 + \sum_{j=0}^{k-1} \beta_j \vec{p}_j, X_k = X_0 + \sum_{j=0}^{k-1} \beta_{kj} \vec{p}_j$

$$\Rightarrow \nabla f(X_k)^T (X - X_k) = \sum_{j=0}^{k-1} \nabla f(X_k)^T \vec{p}_j (\beta_j - \beta_{kj})$$

利用 $\nabla f(X_k)^T \vec{p}_j = 0, j = 0, 1, \dots, k-1$

可得 $\nabla f(X_k)^T (X - X_k) = 0$

再利用凸函数一阶充要条件可得结论

要点：共轭方向二次函数有限终止性

共轭方向二次函数有限终止性

条件: $f(X) = 0.5X^TAX + B^TX + c$, A 对称正定

$\vec{p}_0, \vec{p}_1, \dots, \vec{p}_{n-1}$ 为 A 的共轭方向

$X_0 \in R^n$ 是任意的出发点

由下述直线搜索依次确定 X_1, X_2, \dots, X_n

$$f(X_{k+1}) = f(X_k + t_k \vec{p}_k) = \min_{t>0} f(X_k + t\vec{p}_k)$$

$$k = 0, 1, \dots, n-1$$

结论: $f(X_n) = \min_{X \in R^n} f(X)$

理由：1) 由推论可知

$$f(X_n) = \min \left\{ f(X) \mid \text{s.t. } X = X_0 + \sum_{j=0}^{n-1} \vec{p}_j \beta_j \right\}$$

2) 由原理之一可知 $R^n = \left\{ X \mid X = X_0 + \sum_{j=0}^{n-1} \vec{p}_j \beta_j \right\}$

理由：从共轭方向的几个特点出发：

1、共轭方向线性无关 $\Rightarrow \vec{p}_0, \vec{p}_1, \dots, \vec{p}_{n-1}$ 是 R^n 的一组基

2、 $\nabla f(X) = AX + B$ 是 R^n 的列向量，则对于任意的 \hat{X}

$$\text{若 } 0 \neq \nabla f(\hat{X}) = \alpha_0 \vec{p}_0 + \alpha_1 \vec{p}_1 + \dots + \alpha_{n-1} \vec{p}_{n-1}$$

3、从 X_i 出发沿 \vec{p}_i 直线搜索，则有 $\nabla^T f(X_k + t_k \vec{p}_k) \vec{p}_k = 0$

$$\text{即 } \nabla^T f(X_{k+1}) \vec{p}_k = 0$$

4、 $\nabla f(X_{k+1}) = AX_{k+1} + B = A(X_k + t_k \vec{p}_k) + B = \nabla f(X_k) + t_k A \vec{p}_k$ ，

则有 $\vec{p}_{k-1}^T \nabla f(X_{k+1}) = \vec{p}_{k-1}^T \nabla f(X_k) + \vec{p}_{k-1}^T t_k A \vec{p}_k$ ，进而由3

和共轭方向性质有 $\vec{p}_{k-1}^T \nabla f(X_{k+1}) = 0$ ，依此类推得到

$$\vec{p}_i^T \nabla f(X_{k+1}) = 0, i = 0, 1, \dots, k$$

如果 $\nabla f(X_n) \neq 0$ ，则引发如下矛盾

$$\nabla^T f(X_n) \nabla f(X_n) = \nabla^T f(X_n) (\alpha_0 \vec{p}_0 + \alpha_1 \vec{p}_1 + \dots + \alpha_{n-1} \vec{p}_{n-1}) = 0$$

要点：共轭方向的生成

共轭方向的生成

用Gram-Schmidt 正交化方法顺序生成 A 共轭方向

利用 A 共轭性确定下面方程组中所有待定系数

$$\vec{p}_0 = -\nabla f(X_0)$$

$$\vec{p}_1 = -\nabla f(X_1) + \alpha_{10}\vec{p}_0$$

\vdots

$$\vec{p}_{n-1} = -\nabla f(X_{n-1}) + \alpha_{n-1,0}\vec{p}_0 + \alpha_{n-1,1}\vec{p}_1 + \cdots + \alpha_{n-1,n-2}\vec{p}_{n-2}$$

$$\text{例如: } \vec{p}_0^T A \vec{p}_1 = 0 \Rightarrow 0 = -\vec{p}_0^T A \nabla f(X_1) + \alpha_{10} \vec{p}_0^T A \vec{p}_0$$

$$\Rightarrow \alpha_{10} = \frac{\vec{p}_0^T A \nabla f(X_1)}{\vec{p}_0^T A \vec{p}_0}$$

解前面方程组最终可得

$$\vec{p}_0 = -\nabla f(X_0)$$

$$\vec{p}_1 = -\nabla f(X_1) + \alpha_{10}\vec{p}_0$$

\vdots

$$\vec{p}_{n-1} = -\nabla f(X_{n-1}) + \alpha_{n-10}\vec{p}_0 + \alpha_{n-11}\vec{p}_1 + \cdots + \alpha_{n-1n-2}\vec{p}_{n-2}$$

其中

$$\alpha_{kj} = \frac{\nabla^T f(X_k) \vec{p}_j^T}{\vec{p}_j^T \vec{p}_j}, \quad 1 \leq k \leq n-1, \quad 0 \leq j \leq k-1$$

为了应用于一般性的非线性函数，需要消除 A

消除 A 的基本途经：

$$X_k = X_{k-1} + t_{k-1} \vec{p}_{k-1} \Rightarrow \nabla f(X_k) = \nabla f(X_{k-1}) + t_{k-1} A \vec{p}_{k-1}$$

由推论， $\vec{p}_{k-1}^T \nabla f(X_k) = \vec{p}_{k-1}^T \nabla f(X_{k-1}) + t_{k-1} \vec{p}_{k-1}^T A \vec{p}_{k-1} = 0$

$$\begin{aligned} t_j &= -\frac{\vec{p}_j^T \nabla f(X_j)}{\vec{p}_j^T A \vec{p}_j} & \alpha_{kj} &= \frac{\nabla^T f(X_k) A \vec{p}_j^T}{\vec{p}_j^T A \vec{p}_j} = \frac{t_j \nabla^T f(X_k) A \vec{p}_j^T}{-\vec{p}_j^T \nabla f(X_j)} \\ &\Rightarrow & &= \frac{\nabla^T f(X_k) (\nabla f(X_{j+1}) - \nabla f(X_j))}{-\vec{p}_j^T \nabla f(X_j)} \end{aligned}$$

由于 $j \leq k-1$ ，上式已经可以应用于一般性函数，再利用梯度和共轭方向的关系，可进一步简化系数表达式

$$\vec{p}_j = -\nabla f(X_j) + \alpha_{j0}\vec{p}_0 + \alpha_{j1}\vec{p}_1 + \cdots + \alpha_{jj-1}\vec{p}_{j-1}$$

$$\Rightarrow \nabla f(X_j) = -\vec{p}_j + \alpha_{j0}\vec{p}_0 + \alpha_{j1}\vec{p}_1 + \cdots + \alpha_{jj-1}\vec{p}_{j-1}$$

$$\nabla^T f(X_k) \vec{p}_j = 0, \forall 0 \leq j < k$$

$$\Rightarrow \nabla^T f(X_k) \nabla f(X_j)$$

$$= \nabla^T f(X_k) \left(-\vec{p}_j + \alpha_{j0}\vec{p}_0 + \alpha_{j1}\vec{p}_1 + \cdots + \alpha_{jj-1}\vec{p}_{j-1} \right)$$

$$= 0, \forall 0 \leq j < k$$

$$\nabla^T f(X_k) \nabla f(X_k)$$

$$= \nabla^T f(X_k) \left(-\vec{p}_k + \alpha_{k0}\vec{p}_0 + \alpha_{k1}\vec{p}_1 + \cdots + \alpha_{kk-1}\vec{p}_{k-1} \right)$$

$$= -\nabla^T f(X_k) \vec{p}_k$$

$$\nabla^T f(X_k) \nabla f(X_j) = 0, \quad \forall 0 \leq j < k$$

$$\nabla^T f(X_k) \nabla f(X_k) = -\nabla^T f(X_k) \vec{p}_k$$

$$\begin{aligned} \alpha_{kj} &= \frac{\nabla^T f(X_k) (\nabla f(X_{j+1}) - \nabla f(X_j))}{-\vec{p}_j^T \nabla f(X_j)} \quad \Rightarrow \\ &= \frac{\nabla^T f(X_k) \nabla f(X_{j+1}) - \nabla^T f(X_k) \nabla f(X_j)}{-\vec{p}_j^T \nabla f(X_j)} \end{aligned}$$

$$\alpha_{kj} = \begin{cases} 0 & \text{if } j < k-1 \\ \frac{\nabla^T f(X_k) \nabla f(X_k)}{\nabla^T f(X_{k-1}) \nabla f(X_{k-1})} & \text{if } j = k-1 \end{cases}$$

$$\begin{aligned}
\alpha_{kk-1} &= \frac{\nabla^T f(X_k)(\nabla f(X_k) - \nabla f(X_{k-1}))}{-\vec{p}_{k-1}^T \nabla f(X_{k-1})} \\
&= \frac{\nabla^T f(X_k) \nabla f(X_k)}{\nabla^T f(X_{k-1}) \nabla f(X_{k-1})} && -\vec{p}_{k-1}^T \nabla f(X_{k-1}) \\
& && = 0 - \vec{p}_{k-1}^T \nabla f(X_{k-1}) \\
&= \frac{\|\nabla f(X_k)\|^2}{\|\nabla f(X_{k-1})\|^2} && = \vec{p}_{k-1}^T \nabla f(X_k) - \vec{p}_{k-1}^T \nabla f(X_{k-1}) \\
&= \frac{\nabla^T f(X_k)(\nabla f(X_k) - \nabla f(X_{k-1}))}{\|\nabla f(X_{k-1})\|^2} \\
&= \frac{\nabla^T f(X_k)(\nabla f(X_k) - \nabla f(X_{k-1}))}{\vec{p}_{k-1}^T (\nabla f(X_k) - \nabla f(X_{k-1}))}
\end{aligned}$$

要点：三种共轭梯度法

共轭梯度法 (Fletcher-Reeves)

- 1) 任取 $X_0 \in R^n$, 令 $k = 0$
- 2) 如果 $\|\nabla f(X_k)\| \leq \varepsilon$, 停止计算
- 3) 如果 k/n 等于 0 或整数, 令 $D_k = -\nabla f(X_k)$

否则令 $D_k = -\nabla f(X_k) + \alpha_{k-1}D_{k-1}$, 其中

$$\alpha_{k-1} = \frac{\|\nabla f(X_k)\|^2}{\|\nabla f(X_{k-1})\|^2}$$

- 4) 进行精确搜索获得 $X_{k+1} = X_k + t_k D_k$
- 5) 用 $k + 1$ 替换 k , 回到2) 继续迭代

共轭梯度法 (Polak-Ribiere / Polyak)

- 1) 任取 $X_0 \in R^n$, 令 $k = 0$
- 2) 如果 $\|\nabla f(X_k)\| \leq \varepsilon$, 停止计算
- 3) 如果 k/n 等于 0 或整数, 令 $D_k = -\nabla f(X_k)$

否则令 $D_k = -\nabla f(X_k) + \alpha_{k-1}D_{k-1}$, 其中

$$\alpha_{k-1} = \frac{\nabla^T f(X_k)(\nabla f(X_k) - \nabla f(X_{k-1}))}{\|\nabla f(X_{k-1})\|^2}$$

- 4) 进行精确搜索获得 $X_{k+1} = X_k + t_k D_k$
- 5) 用 $k+1$ 替换 k , 回到 2) 继续迭代

共轭梯度法 (Beale-Sorenson / Hestenes-Stiefel)

- 1) 任取 $X_0 \in R^n$, 令 $k = 0$
- 2) 如果 $\|\nabla f(X_k)\| \leq \varepsilon$, 停止计算
- 3) 如果 k/n 等于 0 或整数, 令 $D_k = -\nabla f(X_k)$

否则令 $D_k = -\nabla f(X_k) + \alpha_{k-1} D_{k-1}$, 其中

$$\alpha_{k-1} = \frac{\nabla^T f(X_k)(\nabla f(X_k) - \nabla f(X_{k-1}))}{D_{k-1}^T (\nabla f(X_k) - \nabla f(X_{k-1}))}$$

- 4) 进行精确搜索获得 $X_{k+1} = X_k + t_k D_k$
- 5) 用 $k+1$ 替换 k , 回到 2) 继续迭代

关于共轭梯度法的结论

- 1) 共轭梯度法是下降算法
- 2) 对于正定二次目标函数

$$f(X) = \frac{1}{2} X^T A X + B^T X + C$$

如果从相同的初始点出发，三种共轭梯度法前进的轨迹完全相同，即每一步一维搜索得到的点均相同，并且，经过 n 次精确的一维搜索后一定找到最优解，即 $X_n = -A^{-1}B$

意义：对一般非线性函数在最优解附近快速收敛

要点：几种算法的性能比较

三种基于梯度的搜索方向的比较

	计算量	效率		鲁棒性
		解附近	远离解	
负梯度	A	C	A	A
共轭梯度	B	B	B	B
牛顿方向	C	A	C	C