

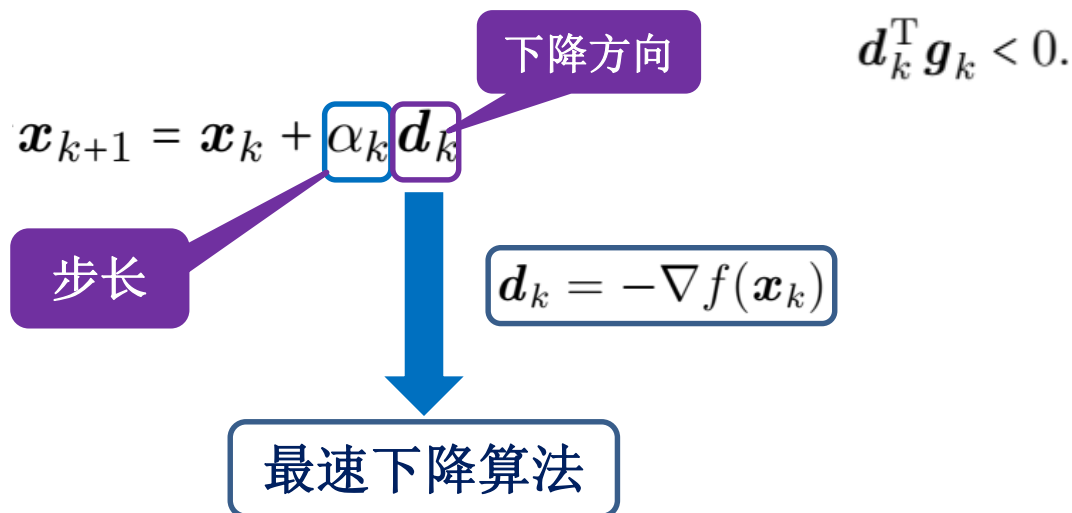
最速下降法与牛顿方法

一、最速下降法

优化模型

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

迭代算法



最速下降算法

步1、取初始点 \mathbf{x}_0 , 参数 $\varepsilon \geq 0$, 令 $k = 0$.

步2、若 $\|\mathbf{g}_k\| \leq \varepsilon$, 算法停止; 否则, 进入下一步.

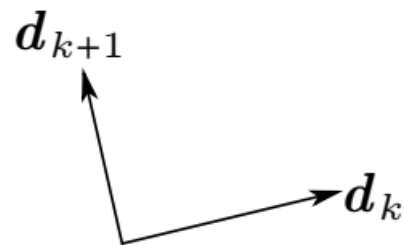
步3、取搜索方向 $\mathbf{d}_k = -\mathbf{g}_k$, 利用精确或非精确线搜索产生步长 α_k .

步4、令 $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k$, $k = k + 1$ 返回步2.

最优步长规则下的收敛性质

迭代方向的正交性 $\mathbf{d}_k^\top \mathbf{d}_{k+1} = 0$.

证明： 记 $\varphi_k(\alpha) = f(\mathbf{x}_k + \alpha \mathbf{d}_k)$.



利用最优步长规则，

$$\begin{aligned} 0 &= \varphi'_k(\alpha_k) = \langle \nabla f(\mathbf{x}_k + \alpha_k \mathbf{d}_k), \mathbf{d}_k \rangle \\ &= \langle \nabla f(\mathbf{x}_{k+1}), \mathbf{d}_k \rangle \\ &= -\langle \mathbf{d}_{k+1}, \mathbf{d}_k \rangle. \end{aligned}$$

算法收敛性

定理 设目标函数连续可微。则最优步长规则下的最速下降算法产生的点列的任一聚点满足 $\nabla f(x^*) = 0$

证明： 将 $d_k = -g_k$ 应用于最优步长下降算法的收敛性结论。

收敛速度估计

定理 对严格凸二次函数 $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \mathbf{G} \mathbf{x}$, $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n > 0$ 为矩阵 \mathbf{G} 的特征根。则最优步长最速下降算法线性收敛, 即

$$\frac{f_{k+1} - f(\mathbf{x}^*)}{f_k - f(\mathbf{x}^*)} \leq \left(\frac{\lambda_1 - \lambda_n}{\lambda_1 + \lambda_n} \right)^2, \quad \frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}_k - \mathbf{x}^*\|} \leq \frac{\lambda_1 - \lambda_n}{\lambda_1 + \lambda_n} \sqrt{\frac{\lambda_1}{\lambda_n}}.$$

目标函数值
线性下降

迭代点列
线性收敛

$\mathbf{x}^* = \mathbf{0}$ 为目标函数的最小值点.

- 很强条件下(二次严格凸)最速下降算法的收敛速度估计;
- 线性收敛是一个很慢的收敛速度;
- 有例为证: 该理论估计是准确的, 没有提升的空间。

例 $\min_{\mathbf{x} \in \mathbb{R}^2} f(x_1, x_2) = \frac{1}{3}x_1^2 + \frac{1}{2}x_2^2$ 唯一最优解 $\mathbf{x}^* = (0; 0)$

取初始点 $\mathbf{x}_0 = (3; 2)$

迭代点列 $\mathbf{x}_k = \left(\frac{3}{5^k}; (-1)^k \frac{2}{5^k} \right) \xrightarrow{k \rightarrow \infty} (0; 0)$

全局收敛，收敛速度线性！ $\|\mathbf{x}_k - \mathbf{x}^*\| = \sqrt{13} \left(\frac{1}{5} \right)^k.$

$$\frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}_k - \mathbf{x}^*\|} = \frac{1}{5} < \frac{1}{5} \sqrt{\frac{3}{2}} = \frac{\lambda_1 - \lambda_2}{\lambda_1 + \lambda_2} \sqrt{\frac{\lambda_1}{\lambda_2}}$$

$$\frac{f(\mathbf{x}_{k+1}) - f(\mathbf{x}^*)}{f(\mathbf{x}_k) - f(\mathbf{x}^*)} = \frac{1}{25} = \left(\frac{\lambda_1 - \lambda_2}{\lambda_1 + \lambda_2} \right)^2$$

算法分析

 搜索方向(最速下降)——下降性是最好的


$$\nabla f(\mathbf{x}_k)^\top \mathbf{d}_k \leq \nabla f(\mathbf{x}_k)^\top \mathbf{d}, \quad \forall \mathbf{d} \in \mathbb{R}^n, \|\mathbf{d}\| = \|\mathbf{d}_k\| = 1$$

 迭代步长(精确搜索)——最优的

$$f(\mathbf{x}_k + \alpha_k \mathbf{d}_k) \leq f(\mathbf{x}_k + \alpha \mathbf{d}_k), \quad \forall \alpha \geq 0$$

 理论性质和数值结果均不好, 较强条件下线性收敛!

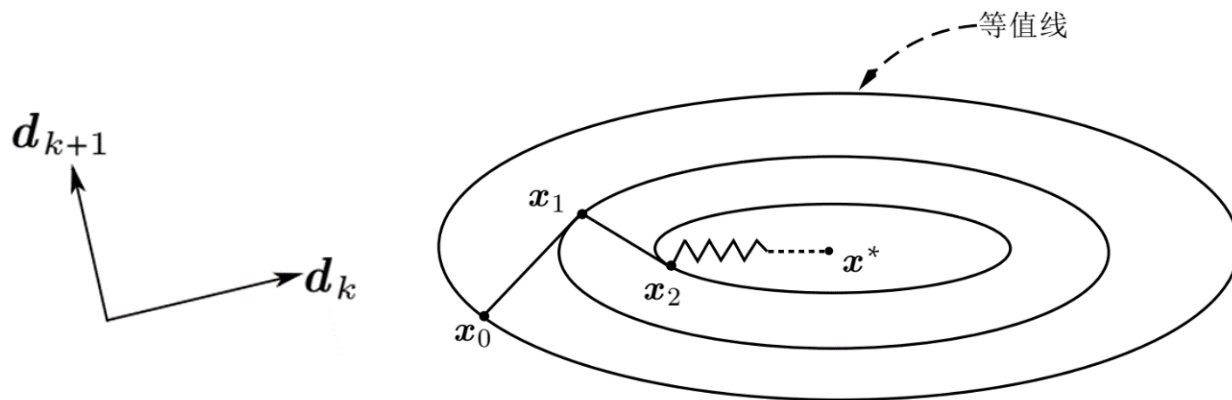
算法分析

 **优点** $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k \quad \mathbf{d}_k = -\mathbf{g}_k$

迭代过程简单，计算量与存储量小；

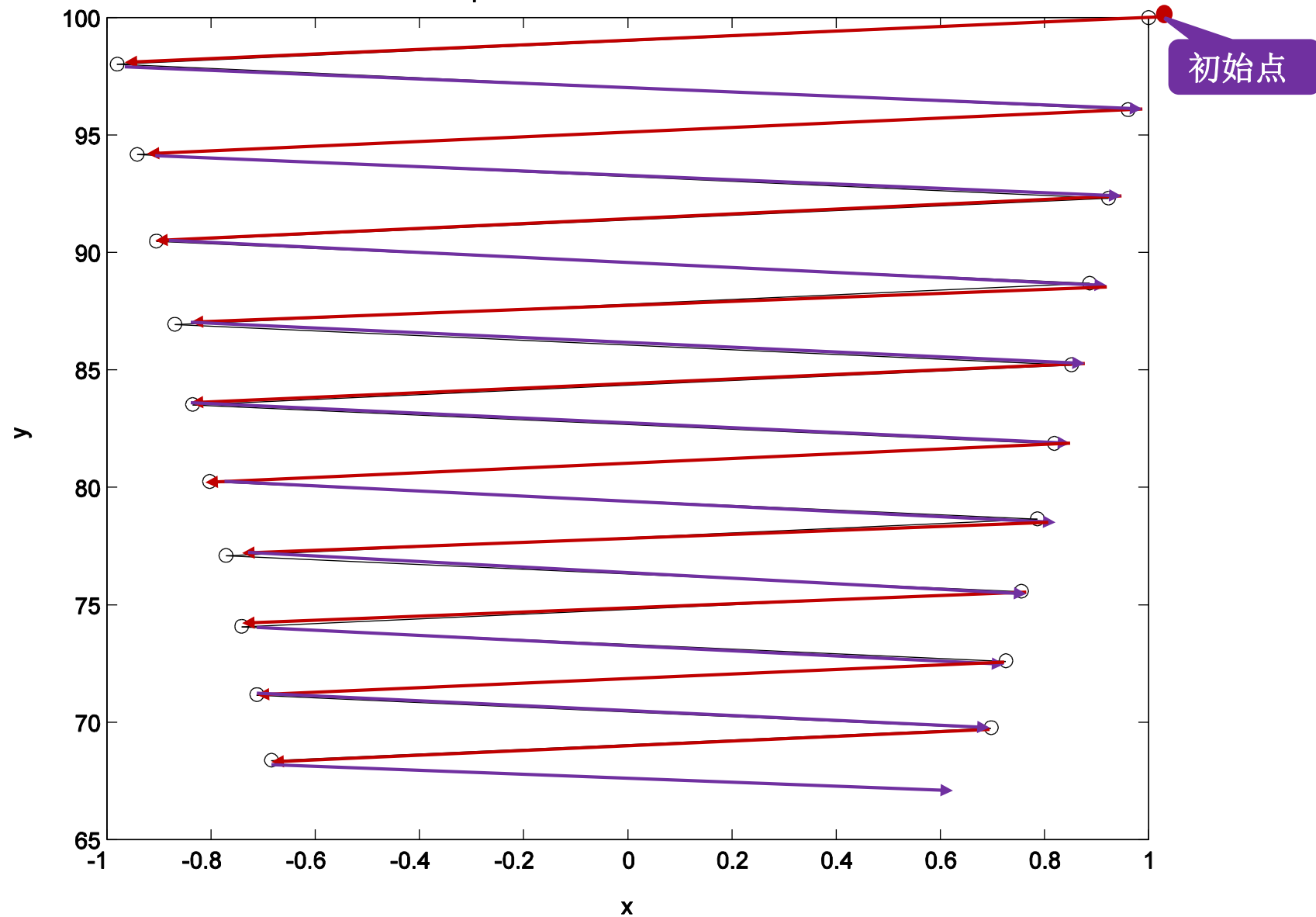
初始步函数值下降较快，可以快速靠近最优解。

 **缺陷** 相邻两搜索方向正交，导致锯齿现象：**越靠近最优值点，靠近速度越慢**。不适于算法收局。



最速下降算法的迭代过程

Steepest descent with exact line search



$$f(\mathbf{x}) = 100x_1^2 + x_2^2$$

原因分析

最速下降法的搜索方向基于目标函数的线性近似

$$f(\mathbf{x}_k + \mathbf{d}) = f(\mathbf{x}_k) + \mathbf{d}^\top \nabla f(\mathbf{x}_k)$$

$$\min\{f(\mathbf{x}_k) + \mathbf{d}^\top \nabla f(\mathbf{x}_k) \mid \|\mathbf{d}\| \leq 1\}$$



二阶近似

牛顿算法

Cauchy-Schwarz不等式

二、牛顿算法

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

利用目标函数 $f(\mathbf{x}_k + \mathbf{d})$ 在 \mathbf{x}_k 点的二阶近似求新的迭代点:

$$m_k(\mathbf{d}) \triangleq f(\mathbf{x}_k) + \mathbf{d}^T \mathbf{g}_k + \frac{1}{2} \mathbf{d}^T \mathbf{G}_k \mathbf{d}.$$

利用最优性条件得其最小值点 $\mathbf{d}_k^N = -\mathbf{G}_k^{-1} \mathbf{g}_k$.

令 $\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{G}_k^{-1} \mathbf{g}_k$ 即得牛顿算法。

牛顿步 $\mathbf{d}_k^N = -\mathbf{G}_k^{-1} \mathbf{g}_k$

算法框架

步1、取初始点 \mathbf{x}_0 , 参数 $\varepsilon \geq 0$, 令 $k = 0$.

步2、若 $\|\mathbf{g}_k\| \leq \varepsilon$, 算法停止; 否则, 进入下一步.

步3、令 $\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{G}_k^{-1} \mathbf{g}_k$, 返回步2.

算法特点: 步长恒取1, 无线搜索;

收敛性与收敛速度

定理 设目标函数二阶连续可微, Hesse阵在最优值点 \mathbf{x}^* 非奇异.

则 (1) 若初始点充分靠近最优值点, 则算法超线性收敛,

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

(2) 若Hesse阵在最优值点附近Lipschitz连续, 则二阶收敛。

证明要点: 利用目标函数二阶展式

$$0 = g(\mathbf{x}^*) = g_k + G_k(\mathbf{x}^* - \mathbf{x}_k) + o(\|\mathbf{x}_k - \mathbf{x}^*\|).$$

左乘 G_k^{-1} , $\mathbf{x}_k - \mathbf{x}^* - G_k^{-1}g_k = o(\|\mathbf{x}_k - \mathbf{x}^*\|)$



$$\mathbf{x}_{k+1} - \mathbf{x}^* = o(\|\mathbf{x}_k - \mathbf{x}^*\|). \text{ 得结论 (1) 。}$$

对后一结论，由目标函数Hesse阵的Lipschitz连续性

$$\begin{aligned}\|\mathbf{x}_{k+1} - \mathbf{x}^*\| &= \|\mathbf{x}_k - \mathbf{x}^* - \mathbf{G}_k^{-1} \mathbf{g}_k\| \\ &= \|\mathbf{G}_k^{-1} [\mathbf{G}_k(\mathbf{x}_k - \mathbf{x}^*) - \mathbf{g}_k + \mathbf{g}(\mathbf{x}^*)]\| \\ &\leq M \left\| \int_0^1 [\mathbf{G}_k - \mathbf{G}(\mathbf{x}^* + \tau(\mathbf{x}_k - \mathbf{x}^*))](\mathbf{x}_k - \mathbf{x}^*) d\tau \right\| \\ &\leq LM \|\mathbf{x}_k - \mathbf{x}^*\|^2 \int_0^1 (1 - \tau) d\tau \\ &= \frac{1}{2} LM \|\mathbf{x}_k - \mathbf{x}^*\|^2.\end{aligned}$$



算法特点

无全局收敛性。仅局部收敛，即要求初始点靠近最优值点；

二阶收敛：迭代点越靠近最优值点，收敛速度越快！

具有二次终止性：对严格凸二次函数，一步即得最优解

严格凸二次函数 $f(x) = \frac{1}{2}x^\top Gx + g^\top x + c$  $x^* = -G^{-1}g$

任意初始点 $x_0 \in R^n$

$$\begin{aligned}x_1 &= x_0 - G^{-1} \nabla f(x_0) \\&= x_0 - G^{-1}(Gx_0 + g) \\&= -G^{-1}g = x^*.\end{aligned}$$



例：用牛顿算法求解 $\min f(x) = \sqrt{1+x^2}$

解：0为最优解。

目标函数导数 $f'(x) = \frac{x}{\sqrt{1+x^2}} \quad f''(x) = \frac{1}{(1+x^2)^{3/2}}$

迭代过程
$$\begin{aligned} x_{k+1} &= x_k - \frac{f'(x_k)}{f''(x_k)} \\ &= x_k - x_k(1+x_k^2) \\ &= -x_k^3 \end{aligned}$$

当 $|x_0| < 1$ 时，算法快速收敛到最优解；

当 $|x_0| \geq 1$ 时，算法不收敛。

例：用牛顿算法求解 $\min f(\mathbf{x}) = 4x_1^2 + x_2^2 - x_1^2x_2$

解： **最优解** $x^* = (0; 0)$ **鞍点** $(2\sqrt{2}; 4)$ $(-2\sqrt{2}; 4)$

目标函数梯度信息

$$\nabla f(\mathbf{x}) = \begin{pmatrix} 8x_1 - 2x_1x_2 \\ 2x_2 - x_1^2 \end{pmatrix} \quad \nabla^2 f(\mathbf{x}) = \begin{pmatrix} 8 - 2x_2 & -2x_1 \\ -2x_1 & 2 \end{pmatrix}$$

取精度 $\varepsilon = 10^{-3}$ 和不同初始点用牛顿算法求解

$$\min f(\mathbf{x}) = 4x_1^2 + x_2^2 - x_1^2x_2$$

初始点 $x_0 = (1; 1)$

目标函数
值非单调

k	$x^{(k)}$	$f(x^{(k)})$	$\nabla f(x^{(k)})$	$\ \nabla f(x^{(k)})\ $	$\nabla^2 f(x^{(k)})$
1	1.0000 1.0000	4.000	6.0000 1.0000	6.0828	6.0000 -2.0000 -2.000 2.0000
2	-0.7500 -1.2500	4.5156	-7.8750 -3.0625	8.4495	10.500 1.5000 1.5000 2.0000
3	-0.1550 -0.1650	0.1273	-1.2911 -0.3540	1.3388	8.3300 0.3100 0.3100 2.0000
4	-0.0057 -0.0111	0.0003	-0.0459 -0.0223	0.0511	8.0222 0.0115 0.0115 2.0000
5	-0.0000 -0.0000	0.0000	-0.0001 -0.0000	0.0001	8.0000 0.0000 0.0000 2.0000

最优解

初始点 $x_0 = (3; 4)$

k	$x^{(k)}$	$f(x^{(k)})$	$\nabla f(x^{(k)})$	$\ \nabla f(x^{(k)})\ $	$\nabla^2 f(x^{(k)})$
1	3.0000	16.000	0.0000	1.0000	0.0000 -6.0000
	4.0000		-1.0000		-6.000 2.0000
2	2.8333	16.0000	0.0000	0.0278	0.0000 -5.6667
	4.0000		-0.2078		-5.6667 2.0000
3	2.8284	16.0000	0.0000	0.0000	0.0000 -5.6569
	4.0000		0.0000		-5.6569 2.0000

鞍点

初始点 $x_0 = (2; 0)$

$$\nabla^2 f(\mathbf{x}_0) = \begin{pmatrix} 8 & -4 \\ -4 & 2 \end{pmatrix}$$

目标函数Hesse阵奇异，不能进行下一步迭代。算法终止

牛顿方法的缺陷

- 初始点远离最优解时，算法可能不收敛；
- Hesse矩阵奇异时，算法不可行；
- 计算量和存储量大：需计算目标函数的梯度和Hesse阵

一些补救和改进措施:

引入线搜索, 阻尼牛顿法, 保证收敛性

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \mathbf{G}_k^{-1} \mathbf{g}_k \quad \longrightarrow \quad \mathbf{x}_{k+1} = \mathbf{x}_k - \boxed{\alpha_k} \mathbf{G}_k^{-1} \mathbf{g}_k$$
$$f(\mathbf{x}_{k+1}) < f(\mathbf{x}_k)$$

结合最速下降法, 设计“杂交”牛顿法, 保证下降性

$$\mathbf{x}_{k+1} = \mathbf{x}_k \boxed{- \mathbf{G}_k^{-1} \mathbf{g}_k} \quad f(\mathbf{x}_{k+1}) \text{ ? } f(\mathbf{x}_k)$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k \boxed{- \alpha_k \mathbf{g}_k}$$

求牛顿方程近似解, 建立非精确牛顿算法, 降低计算量

$$\mathbf{G}_k \mathbf{d} = -\mathbf{g}_k \quad \longrightarrow \quad \mathbf{G}_k \mathbf{d} \boxed{\approx} -\mathbf{g}_k$$

$$\mathbf{G}_k \mathbf{d} = -\mathbf{g}_k + \mathbf{r}_k, \quad \mathbf{r}_k \approx \mathbf{0}$$

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

牛顿算法：求解梯度方程 $g(x) = 0$

最终得到的是目标函数的稳定点： $\nabla f(x^*) = 0$

算法产生点列的聚点，可能是目标函数的最小值点，也可能是最大值点.

将其用于求解恰定方程组 $F(x) = 0$ 得方程组问题的牛顿算法。最常见和最基本的方程组求解算法！

$$\left\{ \begin{array}{l} f_1(x_1, x_2, \cdots, x_n) = 0 \\ f_1(x_1, x_2, \cdots, x_n) = 0 \\ \cdots \cdots \cdots \\ f_m(x_1, x_2, \cdots, x_n) = 0 \end{array} \right. \xrightarrow{\text{ }} \left\{ \begin{array}{l} \text{恰定方程组, } m = n \\ \text{超定方程组, } m > n \\ \text{欠定方程组, } m < n \end{array} \right.$$

一元非线性方程 $f(x) = 0$

取初始点 x_0 ，将函数 $f(x)$ 在该点线性Taylor展开：

$$f(x) \approx f(x_0) + f'(x_0)(x - x_0) = 0$$

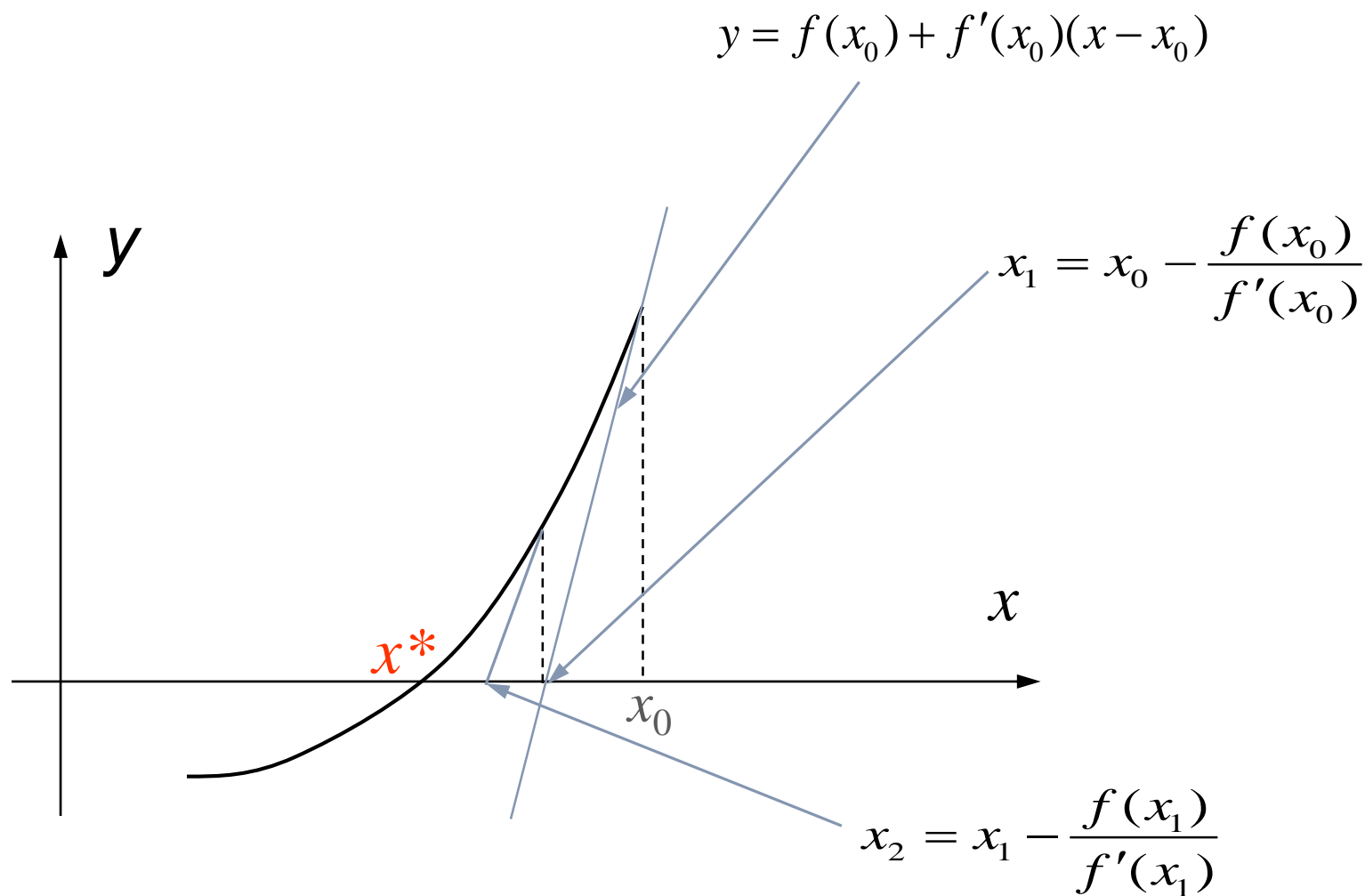
$$x_1 = x_0 - \frac{f(x_0)}{f'(x_0)}$$

重复上述过程

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}$$

定理 一元方程的牛顿迭代算法局部收敛, 在单根附近具有较高(二阶)的收敛速度.

条件： $f'(x^*) \neq 0$



牛顿法，即切线法

$$\begin{cases} f_1(x_1, x_2, \dots, x_n) = 0 \\ f_2(x_1, x_2, \dots, x_n) = 0 \\ \vdots \\ f_n(x_1, x_2, \dots, x_n) = 0 \end{cases} \quad \longrightarrow \quad F(\mathbf{x}) = 0$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \nabla F(\mathbf{x}_k)^{-1} F(\mathbf{x}_k)$$

$$\nabla F(\mathbf{x}) = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \dots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \dots & \frac{\partial f_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_n}{\partial x_1} & \frac{\partial f_n}{\partial x_2} & \dots & \frac{\partial f_n}{\partial x_n} \end{bmatrix}$$

定理 多元恰定方程的牛顿迭代算法局部收敛, 在根附近具有较高(二阶)的收敛速度.

条件: Jacob矩阵非奇异。

$$\nabla F(\mathbf{x}) = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \cdots & \frac{\partial f_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_n}{\partial x_1} & \frac{\partial f_n}{\partial x_2} & \cdots & \frac{\partial f_n}{\partial x_n} \end{bmatrix}$$

$$\begin{cases} a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n = b_1 \\ a_{21}x_1 + a_{22}x_2 + \cdots + a_{2n}x_n = b_2 \\ \vdots \quad \quad \quad \vdots \quad \quad \quad \vdots \quad \quad \quad \vdots \\ a_{n1}x_1 + a_{n2}x_2 + \cdots + a_{nn}x_n = b_n \end{cases}$$

系数矩阵可逆

$$\mathbf{A}\mathbf{x} = \mathbf{b} \longrightarrow \mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$$

两种算法的比较

	最速下降算法	牛顿算法
收敛性	全局收敛 (初始点任意)	局部收敛 (初始点靠近最优值点)



最速下降算法：无论身在何处都能找回家，只是速度慢点，特别到了家门口。

牛顿算法：只能在家附近才能找到家，而且此时能快速回到家，只是累点。