

# 第一章

## 1. Introduction

optimization problem

$$\begin{aligned} \min \quad & f(x) \\ \text{s.t.} \quad & g(x) \leq 0 \\ & h(x) = 0 \end{aligned}$$

$x = (x_1, \dots, x_n) \in \mathbb{R}^n$ : optimization variables

$f: \mathbb{R}^n \mapsto \mathbb{R}$ : objective function

$g: \mathbb{R}^n \mapsto \mathbb{R}^{m_g}$ : inequality constraint functions

$h: \mathbb{R}^n \mapsto \mathbb{R}^{m_h}$ : equality constraint functions

**optimal solution**  $x^*$  has smallest value  $f$  among all vectors that satisfy the constraints.

前提:

1.  $f(x)$  is lower bounded 有下界
2.  $f(s)$  has bounded sub-level set 有有界的水平集

## 2. Convex Sets

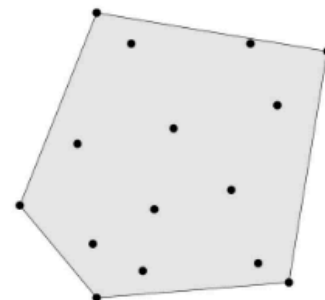
### Convex Sets

Def: A set is convex if every line between its two points stays in the set?

$$\theta x_1 + (1 - \theta)x_2, \quad 0 \leq \theta \leq 1$$

More General:  
All **convex combinations**  
lie in set.

$$\begin{aligned} \sum \theta_i x_i &= \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 \\ \sum \theta_i &= \theta_1 + \theta_2 + \theta_3 = 1, \quad \theta_i \geq 0 \forall i \end{aligned}$$



- 凸组合: 表示凸包
- 重心坐标: Barycentric coordinate

凸集中任意一点的凸表示形式:  $\theta_i$  表示广义重心坐标

- convex hull: 集合中所有点的凸组合, 找一个最小的凸集把非凸的集合包裹起来

### 3. High-Order Info of Functions

- 高阶导数信息

Gradient  $\mathbb{R} \mapsto \mathbb{R}^n$

Hessian  $\mathbb{R} \mapsto \mathbb{R}^{n \times n}$

Hessian是Gradient的Jacobian

光滑函数的Hessian是对称的

**光滑函数**（英语：Smooth function）在**数学**中特指无穷**可导**的函数，不存在**尖点**，也就是说所有的有限**阶**导数都存在。例如，**指数函数**就是光滑的，因为指数函数的导数是指数函数本身。

若一函数是**连续**的，则称其为 $C^0$ 函数；若函数存在导函数，且其导函数连续，则称为**连续可导**，记为 $C^1$ 函数；若一函数 **$n$ 阶可导**，并且其 **$n$ 阶导函数连续**，则为 $C^n$ 函数（ $n \geq 1$ ）。而光滑函数是对所有 **$n$ 都属于 $C^n$ 函数**，特称其为 $C^\infty$ **函数**。

- 

Function  $f(x) = f(x_1, x_2, x_3)$

Gradient  $\nabla f(x) = \begin{pmatrix} \partial_1 f(x) \\ \partial_2 f(x) \\ \partial_3 f(x) \end{pmatrix}$

Hessian  $\nabla^2 f(x) = \begin{pmatrix} \partial_1^2 f(x) & \partial_1 \partial_2 f(x) & \partial_1 \partial_3 f(x) \\ \partial_2 \partial_1 f(x) & \partial_2^2 f(x) & \partial_2 \partial_3 f(x) \\ \partial_3 \partial_1 f(x) & \partial_3 \partial_2 f(x) & \partial_3^2 f(x) \end{pmatrix}$  Symmetric for smooth func

Linear approx

Approx at zero  $f(x) = \underbrace{f(0) + x^T \nabla f(0) + \frac{1}{2} x^T \nabla^2 f(0) x}_{\text{Quadratic approx}} + O(\|x - x_0\|^3)$

Quadratic approx

#### Jacobian

The extension of the gradient of multidimensional  $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}^m$ :

$$f'(x) = \frac{df}{dx^T} = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \cdots & \frac{\partial f_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \frac{\partial f_m}{\partial x_2} & \cdots & \frac{\partial f_m}{\partial x_n} \end{pmatrix}$$

$$f(x) : X \rightarrow Y; \quad \frac{\partial f(x)}{\partial x} \in G$$

X	Y	G	Name
$\mathbb{R}$	$\mathbb{R}$	$\mathbb{R}$	$f'(x)$ (derivative)
$\mathbb{R}^n$	$\mathbb{R}$	$\mathbb{R}^n$	$\frac{\partial f}{\partial x_i}$ (gradient)
$\mathbb{R}^n$	$\mathbb{R}^m$	$\mathbb{R}^{n \times m}$	$\frac{\partial f_i}{\partial x_j}$ (jacobian)
$\mathbb{R}^{m \times n}$	$\mathbb{R}$	$\mathbb{R}^{m \times n}$	$\frac{\partial f}{\partial x_{ij}}$

- 分子分母布局可以确定结果的行列

The notation of differentials could be useful here

$$dA = 0$$

$$d(\alpha X) = \alpha(dX)$$

$$d(AXB) = A(dX)B$$

$$d(X + Y) = dX + dY$$

$$d(X^\top) = (dX)^\top$$

$$d(XY) = (dX)Y + X(dY)$$

$$d\langle X, Y \rangle = \langle dX, Y \rangle + \langle X, dY \rangle$$

$$d\left(\frac{X}{\phi}\right) = \frac{\phi dX - (d\phi)X}{\phi^2}$$

$$d \operatorname{tr} X = I$$

$$df(g(x)) = \frac{df}{dg} \cdot dg(x)$$

## 4. Convex Functions

$f(x)$ 上任意两点的线段一定在函数上方

local minimum 一定是 global minimum, 并且解集一定是凸的

凸函数有凸的次水平集

仿射变换  $Ax + b$  保留凸性

**max**保留凸性 ( $f(x)$ 定义为几个凸函数中的最大函数)

Other allowed operations

Set sum  $A + B = \{x + y \mid x \in A, y \in B\}$

Set product  $A \times B = \{(x, y) \mid x \in A, y \in B\}$

Point-wise max preserves convexity

$$g(x) = \max_i f_i(x)$$

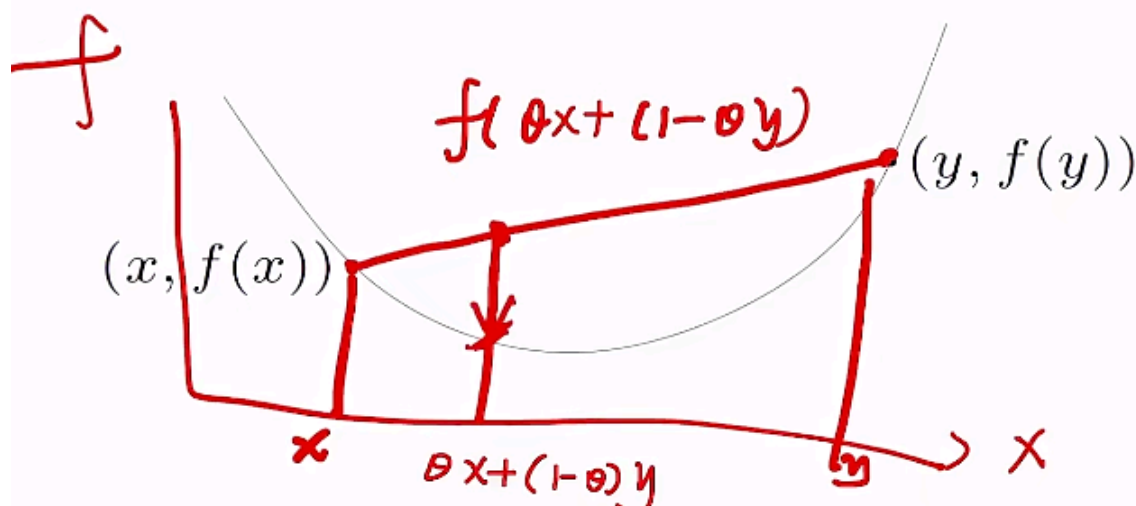
Absolute value  $|x| = \max\{x, -x\}$

Infinity norm  $\|x\|_\infty = \max_i |x_i|$

Max eigenvalue  $\|A\|_2 = \max_v v^T A v$

Jensen's inequality

$$f(\theta x + (1 - \theta)y) \leq \theta f(x) + (1 - \theta)f(y)$$



- 凸函数:

## Convex Functions

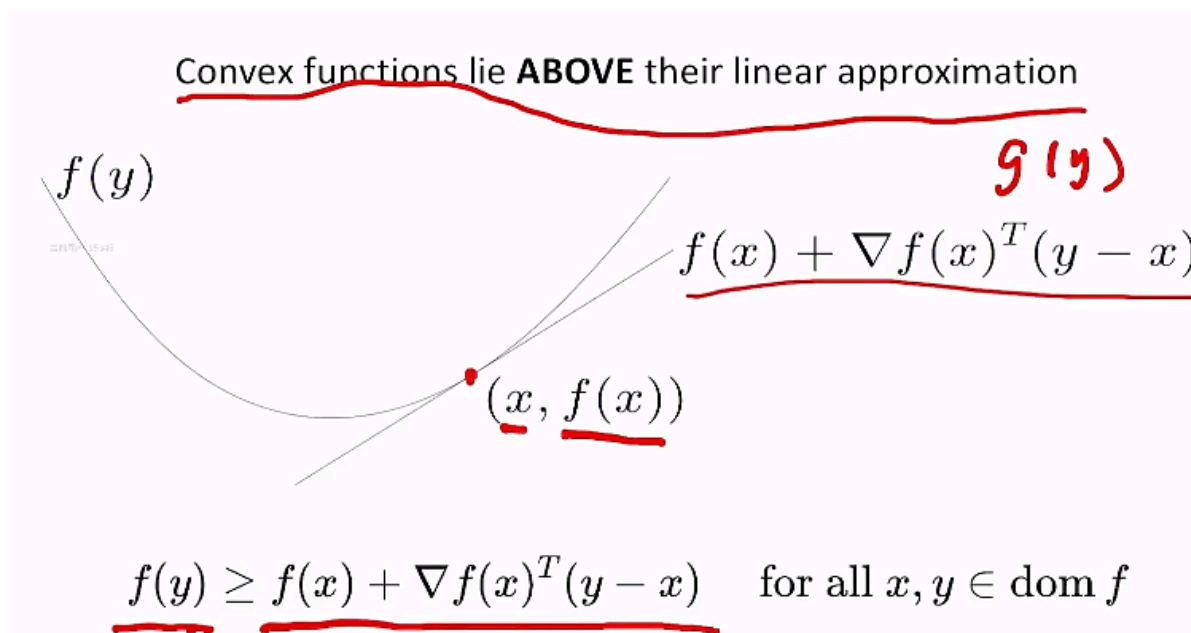
Why are these convex?

Trace	$f(X) = \text{trace}(A^T X)$	Linear operator
Distance over set	$f(x) = \max_{y \in C} \ x - y\ $	Max over convex
Distance to convex set	$f(x) = \min_{y \in C} \ x - y\ $	Special case
Affine Norm	$f(x) = \ b + \sum_i A_i x_i\ _2$	Affine comp

If  $g(x,y)$  is convex, then minimizing for  $y$  preserves convexity

- 凸函数性质

整个函数在某点处切线的上方



梯度为0，全局最优，若函数非凸，则为局部最优

- 若函数光滑

光滑性类别	梯度 $\nabla f(x)$	Hessian $H_f(x)$	混合偏导对称性
$C^0$ (连续)	不一定存在	不一定存在	-
$C^1$ (一阶可微)	存在且连续	可能不存在	可能不对称
$C^2$ (二阶可微)	存在且连续	存在且连续	确保对称
$C^k$ ( $k$ 阶可微)	存在且连续	存在且可微	对称 (只要 $C^2$ )
$C^\infty$ (光滑)	无限阶可微	无限阶可微	对称

1. Hessian 矩阵的存在性要求函数至少是  $C^2$ ，即所有二阶偏导数都存在且连续。
2. 如果函数仅是二阶可微（但二阶导数不连续），Hessian 仍然存在，但混合偏导不一定对称。
3. Hessian 对称性需要  $C^2$  条件（Schwarz 定理）。
4. 光滑函数（ $C^\infty$ ）不仅保证 Hessian 矩阵存在，还能保证所有高阶导数存在。
5. 非光滑函数（如 ReLU）可能没有 Hessian 矩阵，例如：

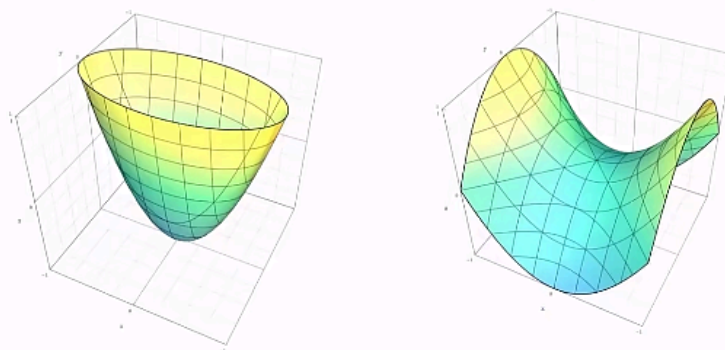
$$f(x) = \max(0, x)$$

在  $x = 0$  处不可导，所以 Hessian 矩阵不存在 ↓

### Second-order conditions

A smooth function is convex iff  $\nabla^2 f(x) \succeq 0, \quad \forall x$  (semi-definite)

For non-convex functions, minima satisfy:  $\nabla^2 f(x^*) \succeq 0$



The Hessian is a good local model of a smooth function



光滑函数或Hessian存在（至少 $C^2$ ）的函数为凸当且仅当任意点处Hessian是半正定

非凸光滑函数**local minima**处的Hessian半正定

#### • 强凸性质

一个凸函数，Hessian严格正定，则函数为强凸的

#### 1. 得到 $f(x)$ 下界

## Strong convexity

$$f(y) \geq \underbrace{f(x) + (y-x)^T \nabla f(x)}_{\text{holds for any convex func}} + \underbrace{\frac{m}{2} \|y-x\|^2}_{\text{min curvature}}$$

When Hessian exists

$$\begin{aligned} f(y) &\approx f(x) + (y-x)^T \nabla f(x) + \frac{1}{2} (y-x)^T \nabla^2 f(x) (y-x) \\ &\geq f(x) + (y-x)^T \nabla f(x) + \frac{\lambda_{\min}}{2} \|y-x\|^2 \end{aligned}$$


and so

$$\nabla^2 f(x) \succeq mI$$

2. 得到f(x)上界

The Lipschitz constant  $M$  of  $\nabla f(x)$  satisfies

$$\|\nabla f(x) - \nabla f(y)\| \leq M \|y-x\|$$


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

3. 约束函数值的界， $m$ 、 $M$ 可以用来判别函数的条件数

We can bound the objective error in terms of distance from minimizer

$$f(y) - f(x^*) \geq \frac{m}{2} \|y - x^*\|^2 \qquad f(y) - f(x^*) \leq \frac{M}{2} \|y - x^*\|^2$$

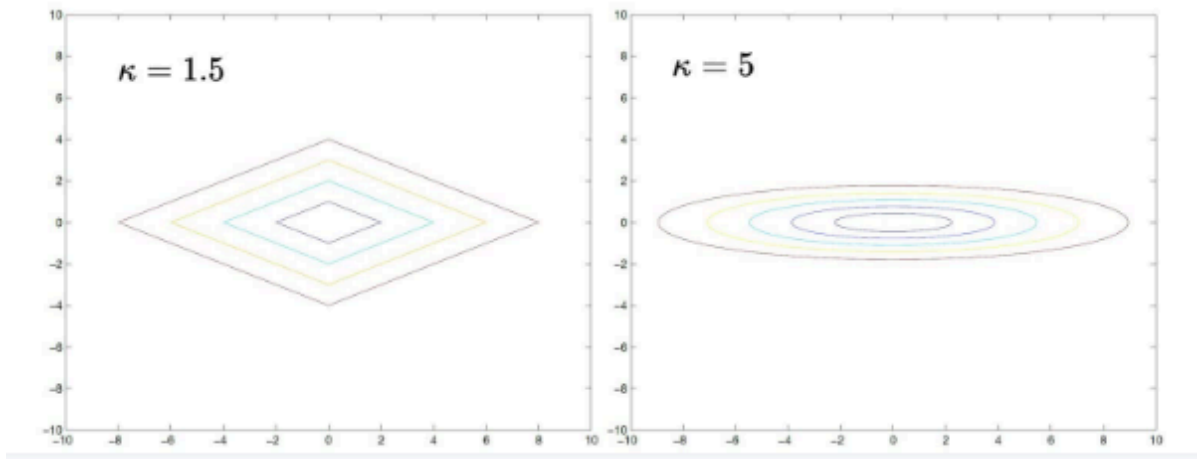
(暂时没懂) 条件数

## Condition number

For any function  $\kappa = \frac{\text{major axis}}{\text{minor axis}}$

For smooth functions  $\kappa \approx \text{cond}(\nabla^2 f(x))$

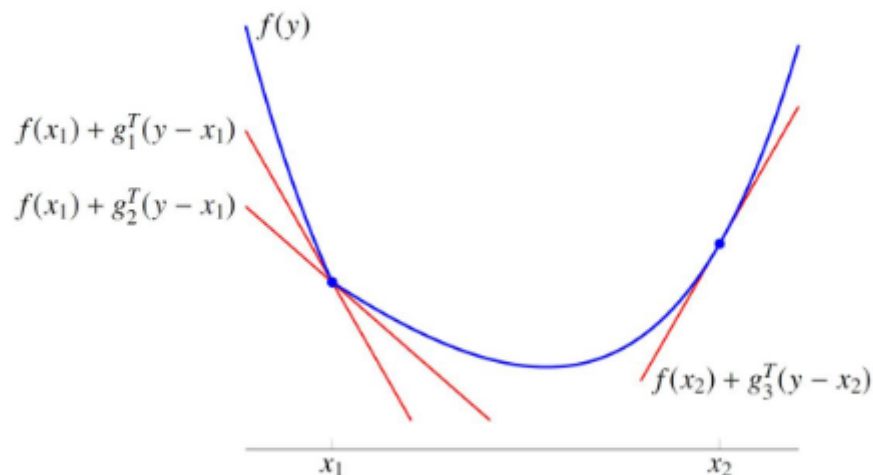
For differentiable functions  $\kappa = M/m$



- 凸但是不可微的函数（不可微的凸函数的全局极小值）

## Sub-differential

$$\partial f(x) = \{g : f(y) \geq f(x) + (y - x)^T g, \forall y\}$$



Optimality:  $0 \in \partial f(x^*)$

1. 可微的点处为切线，连续但不可微的地方为梯度的凸组合（次梯度）
  2. 最优性：0属于次梯度（比如 $f(x)=|x|$ ，0属于 $\{-1,1\}$ 的凸组合）
  3. 不光滑的函数沿着次梯度反方向走函数值可能不下降反而上升（对收敛速度造成挑战）  
需要注意的是负的次梯度方向不一定是函数值下降方向，而只有方向导数  $<0$  的方向才是函数值下降方向
  4. 如果知道non-smooth的分界点时， $f(x)$ 下降最快的方向为non-smooth点处次梯度模长最小向量的反方向（如何证明？），可以加快收敛速，防止振荡
- 凸函数的梯度/次梯度单调性



Monotonicity: The (sub) gradient of any convex func is monotone

$$\langle y - x, \nabla f(y) - \nabla f(x) \rangle \geq 0$$

or

$$\langle y - x, g_y - g_x \rangle \geq 0, g_x \in \partial f(x), g_y \in \partial f(y)$$

This can be obtained by adding two equations below.

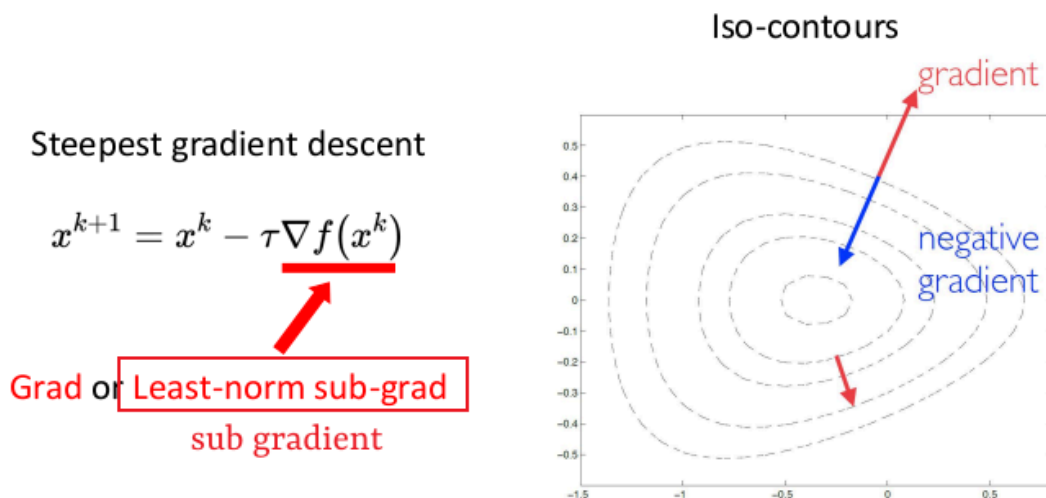
$$f(y) \geq f(x) + \nabla f(x)^T (y - x)$$

$$f(x) \geq f(y) + \nabla f(y)^T (x - y)$$

如果函数严格凸，取严格大于号

## 5. Unconstrained Optimization for Nonconvex Functions

- Steepest Gradient Descent 最速梯度下降



**exact line search**本质上要求 $\min f(x_{k+1})$ ，采用了当前“最好”的方向，和“最好”的步长因子。因为最好的方向就是梯度的反向，最好的步长因子满足

$$\alpha_k = \operatorname{argmin}_{\alpha} f(x_k + \alpha d_k)$$

```

while || ∇ f(xk) || ≥ ε
    dk = - ∇ f(xk)
    αk = argminα f(xk + α dk)
    xk+1 = xk + αk dk
    k = k + 1
end

```

每一步仍要求一个子的优化问题，这个代价一般上比较大，但是对于二次函数可以方便求极值，二次函数，所以才可能直接求出最优的步长，对于任意非线性函数，最优步长一般是无法求得解析解的，所以才会有Backtracking Line Search方法。

最速下降法使用的相邻下降方向是正交的。在最速下降法中，当次迭代的梯度方向也是和上次迭代梯度方向垂直，但和再之前的梯度方向就不垂直了，所以会有“之”形路线。然后，共轭梯度法要求的是关于矩阵正交，并非直接正交。这是因为每一步都将一个方向走到最优，对于n维空间，那么只需走n步，每一步走的方向都是一个维度。（注意每一步走的方向所代表的维度不一定与坐标轴平行（我们这里说的维度不是坐标轴，是该空间的任何基中的一个方向空间的维数是什么）。但是每个方向之间一定要正交，所以表现出来共轭梯度只需要n步）

### 1. 如何确定搜索的步长

- Constant step size  $\tau = c$  Not intelligent
- Diminishing step size  $\tau = c/k$  Robbins-Monro rule for expensive func
- Exact line search  $\tau = \arg \min_{\alpha} f(x^k + \alpha d)$  Generally nontrivial
- Inexact line search  $\tau \in \{\alpha \mid f(x^k) - f(x^k + \alpha d) \geq -c \cdot \alpha d^T \nabla f(x^k)\}$

Armijo condition, easy to satisfy

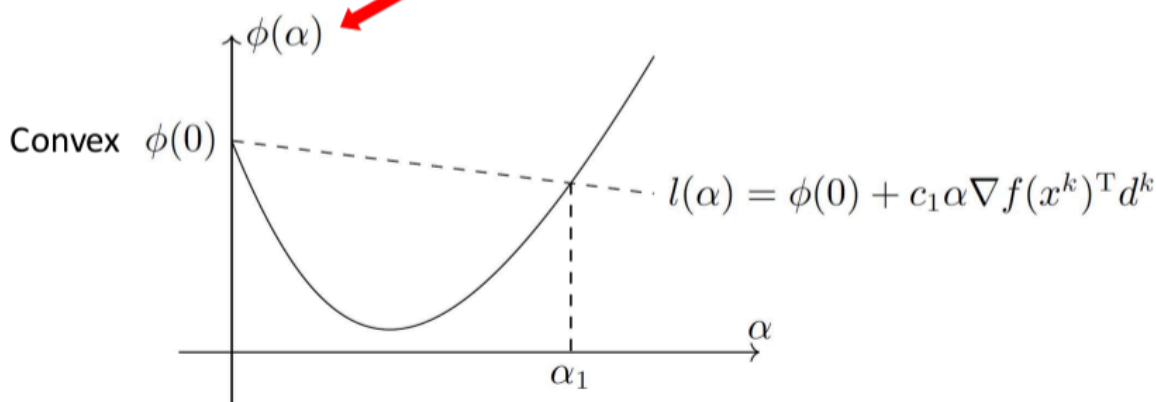
好的步长：使函数充分下降，不能太松弛（收敛很慢）也不能太激进（发散震荡）

### 2. Backtracking/Armijo Line search(本质上是让 $f(x_k + \alpha d)$ 充分下降)

Armijo condition (sufficient decrease condition)

$$\tau \in \{\alpha \mid f(x^k) - f(x^k + \alpha d) \geq -c \cdot \alpha d^T \nabla f(x^k)\}$$

$$c \in (0, 1)$$



实际使用中可以用Armijo condition（充分下降条件）作为选择步长是否合适的判断准则，用二分法找合适的步长（适用于光滑或者分片光滑，凸或者非凸都可以收敛到local minima）

Choose search direction:  $d = -\nabla f(x^k)$

While  $f(x^k + \tau d) > f(x^k) + c \cdot \tau d^T \nabla f(x^k)$

$$\tau \leftarrow \tau/2$$

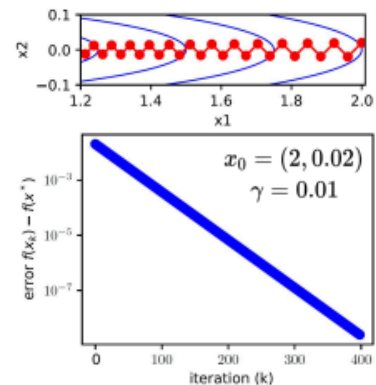
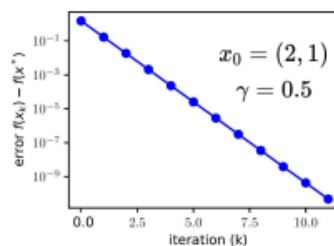
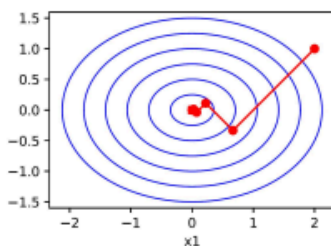
Update iterate  $x^{k+1} = x^k + \tau d$

Repeat this until **gradient is small**  
or **subdifferential contains zero**.

### 3. 最速下降法的缺点

Drawbacks: Poor conditioning causes performance degeneration

$$f(x_1, x_2) = \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} = \frac{\gamma}{2} x_1^2 + \frac{1}{2} x_2^2, \quad \mathbf{Q} = \text{diag}\{\gamma, 1\}$$



**Curvature info is needed!**

对于条件数比较敏感，原因是最速梯度下降只用了函数的一阶信息，而条件数表明了函数的高阶信息，包含了函数的曲率信息，表示函数被压的多宽多扁，不利用曲率信息就会不断震荡

## 6. Modified Damped Newton's Method

要求:  $f(x), f'(x), f''(x)$  连续, Hessian 连续, 函数可以是非凸

- Newton's Method

## Newton's Method

By second-order Taylor expansion,

$$f(\mathbf{x}) \approx \hat{f}(\mathbf{x}) \triangleq f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T (\mathbf{x} - \mathbf{x}_k) + \frac{1}{2} (\mathbf{x} - \mathbf{x}_k)^T \nabla^2 f(\mathbf{x}_k) (\mathbf{x} - \mathbf{x}_k)$$

Minimizing quadratic approximation

$$\nabla \hat{f}(\mathbf{x}) = \nabla^2 f(\mathbf{x}_k) (\mathbf{x} - \mathbf{x}_k) + \nabla f(\mathbf{x}_k) = \mathbf{0}$$

$$\Rightarrow \mathbf{x} = \mathbf{x}_k - [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k)$$

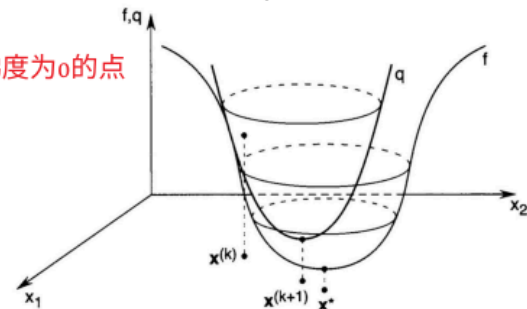
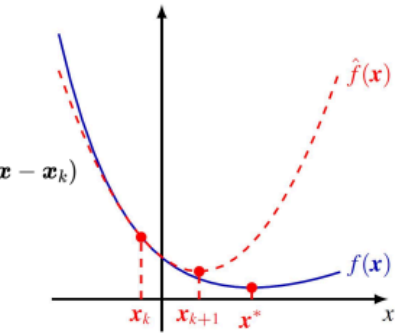
provided  $\nabla^2 f(\mathbf{x}_k) \succ \mathbf{0}$  Hessian 严格正定 充要

Newton step

函数严格凸 可以直接求梯度为0的点

$$\mathbf{x}_{k+1} = \mathbf{x}_k - [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k)$$

Note. If the function is quadratic, then Newton's method gets to the optimum in a single step.

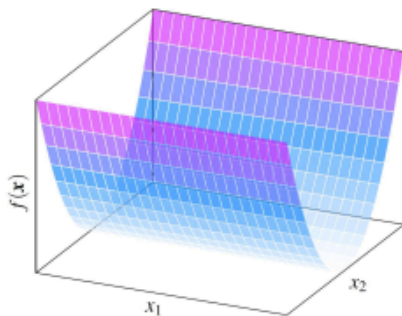


若f(x)本身是严格正定的二次函数，牛顿法可以一步收敛到精确解

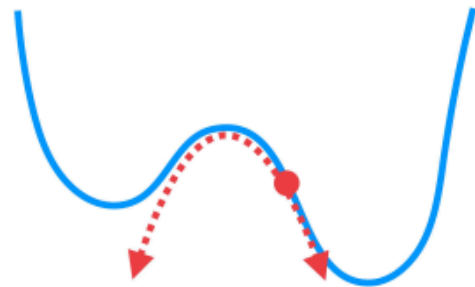
缺点：每次迭代需要求一次Hessian矩阵的逆，Hessian必须严格正定，需要让牛顿步长与负梯度方向夹角小于90度

Drawbacks: In practice Hessian can be singular and indefinite

$$\mathbf{x}_{k+1} = \mathbf{x}_k - [\nabla^2 f(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k)$$



Singular Hessian!



If start at a point of negative curvature, the Newton goes up!

(negative) search direction: must form an acute angle with the gradient

评判数值优化算法的标准：

Three aspect to evaluate a numerical optimization method:

1. Convergence speed (How to measure the rate? In Lec2.).
2. Stability when applied to different functions.
3. Computation work per iteration.

- Practical Newton's Method

```
initialization  $\mathbf{x} \leftarrow \mathbf{x}_0 \in \mathbb{R}^n$ 
while  $\|\nabla f(\mathbf{x})\| > \delta$  do
     $\mathbf{d} \leftarrow -\mathbf{M}^{-1}\nabla f(\mathbf{x})$ 
     $t \leftarrow$  backtracking line search
     $\mathbf{x} \leftarrow \mathbf{x} + t\mathbf{d}$ 
end while
return
```

- Choose a positive-definite  $\mathbf{M}$  that is close to the Hessian 用严格正定M代替Hessian
- Solve the linear system via factorization instead of inversion 用linear solver代替矩阵求逆
- line search does not need grad and Hessian line search 只需要函数值的信息  
不需要gradient和hessian

#### 1. Hessian半正定

- 函数是凸的

- If function is convex, its Hessian must be PSD. We choose a  $\mathbf{M}$  as

$$\mathbf{M} = \nabla^2 f(\mathbf{x}) + \epsilon \mathbf{I}, \quad \epsilon = \min(1, \|\nabla f(\mathbf{x})\|_\infty)/10$$

Since  $\mathbf{M}$  is PD, the search direction is solved by Cholesky factorization

$$\mathbf{M}\mathbf{d} = -\nabla f(\mathbf{x}), \quad \mathbf{M} = \mathbf{L}\mathbf{L}^T$$

where  $\mathbf{L}$  is a lower triangular matrix.

- 函数是非凸的

- If function is nonconvex, its Hessian can be indefinite. We compute  $\mathbf{M}$  on the fly through

Bunch-Kaufman Factorization

$$\mathbf{M}\mathbf{d} = -\nabla f(\mathbf{x}), \quad \mathbf{M} = \mathbf{L}\mathbf{B}\mathbf{L}^T$$

where  $\mathbf{B}$  is block diagonal matrix with block size 1x1 or 2x2. All 2x2 blocks contain nonpositive eigenvalues, which are easy to modify.

把所有2x2矩阵负的特征值都变为 $\epsilon$ 的极小量

稀疏矩阵Hessian分解可以进一步加快