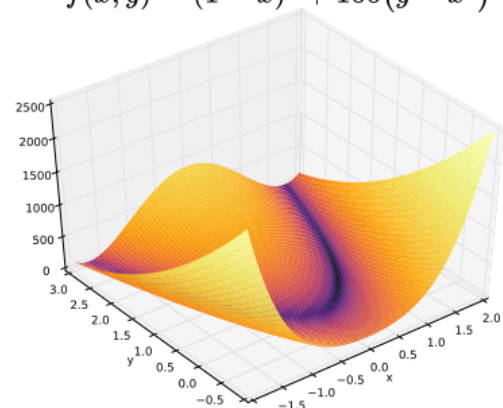


第二章

1. Why Quasi-Newton Method

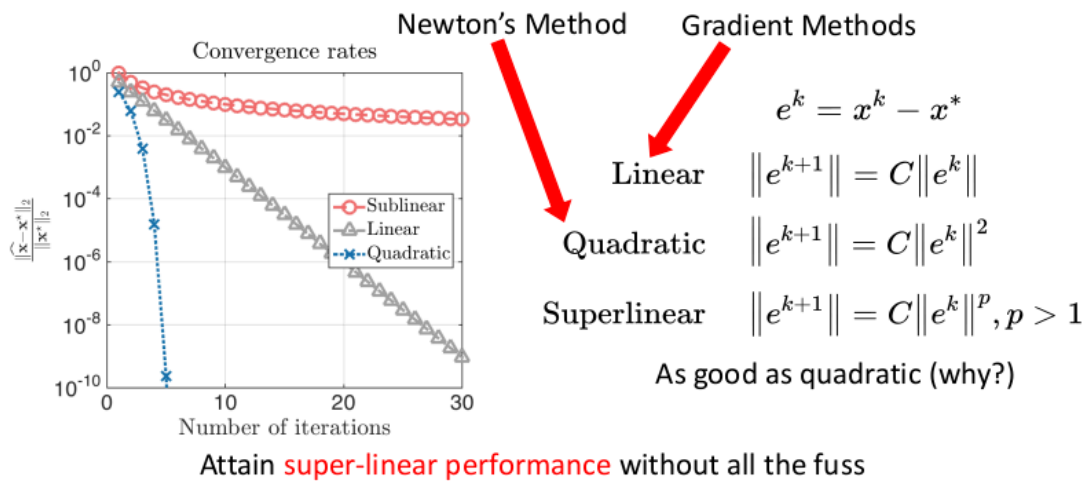
Why Quasi-Newton Methods

Rosenbrock banana function
 $f(x, y) = (1 - x)^2 + 100(y - x^2)^2$



- Gradient descent is cheap
 - Compute partials: $O(N)$
- Newton's method is expensive
 - Compute mixed partials: $O(N^2)$
 - Linear solver: $O(N^3)$
- Other problems
 - When quadratic approximation is bad, Newton's is a waste
 - Hessian becomes poorly conditioned
 - Nonconvex problems: indefinite Hessian = ascent step
- 梯度法
 - 优点：计算梯度的时间复杂度为 $O(N)$
 - 缺点：没有函数的曲率信息，可能收敛慢
- 牛顿方法
 - 优点：包含曲率信息，收敛快
 - 缺点：计算Hessian复杂度为 $O(N^2)$ ，计算最优解设计Hessian求逆，复杂度为 $O(N^3)$
- 其他问题：
 - 非线性函数使用Newton方法近似效果很差的话，会造成计算Hessian和求逆资源浪费
 - Hessian奇异的时候牛顿方法不稳定
 - 非凸问题Hessian有负特征值， $f(x)$ 可能会上升到maxmia
- 期望以较低的复杂度实现超线性的收敛速率 **linear < superlinear < quadratic**

Normally, for **convex smooth** functions, we have



2. Quasi-Newton Methods (凸且光滑函数)

Newton approximation (PD Hessian H)

$$f(x) - f(x^k) \approx (x - x^k)^T g + \frac{1}{2} (x - x^k)^T H (x - x^k)$$

$$\text{Solve } H^k d^k = -g^k$$

Quasi-Newton approximation

$$f(x) - f(x^k) \approx (x - x^k)^T g^k + \frac{1}{2} (x - x^k)^T M^k (x - x^k)$$

$$\text{Solve } M^k d^k = -g^k$$

What conditions should M satisfy?

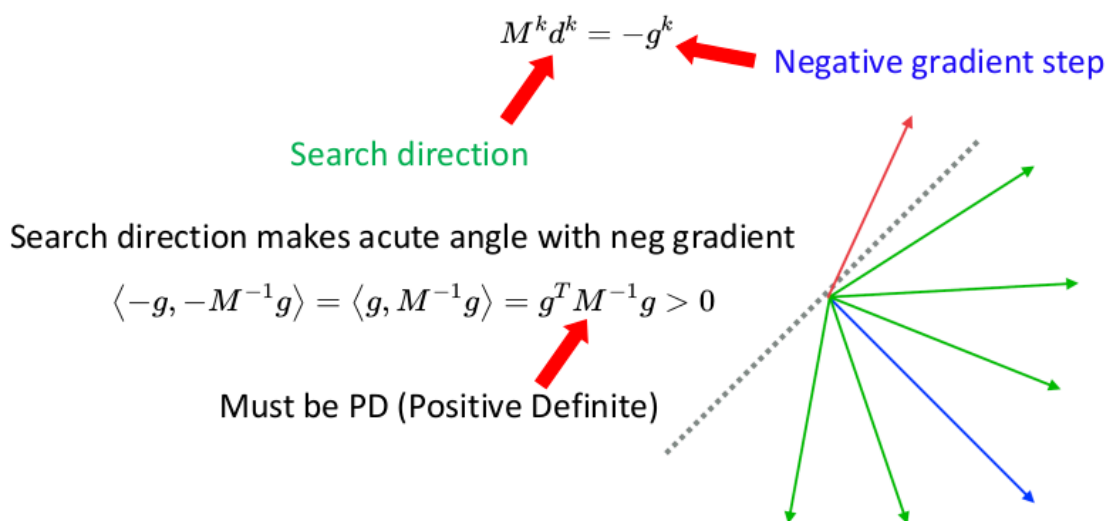
- It should not need full 2nd-order derivatives
- The linear equations have closed-form solutions
- It should be lightweight and compact to store
- **It must preserve descent directions**
- **It should contain curvature information (local quadratic approx.)**

1. M 不需要所有二阶导信息
2. M 需要构造成容易解线性方程组的形式，有闭式解
3. M 必须是稀疏 轻量 紧凑
4. 方向 d 必须是函数值下降的方向（ M 严格正定来保证）
5. 必须包含函数曲率信息，Hessian的局部近似（）

拟牛顿法的想法其实很简单，就像是函数值的两点之差可以逼近导数一样，一阶导数的两点之差也可以逼近二阶导数。几何意义是求一阶导数的“割线”，取极限时，割线会逼近切线，矩阵 B 就会逼近Hessian矩阵的逆。

（ M 严格正定来保证函数值下降方向）

When is descent direction preserved?



(不必求Hessian和求逆操作，用梯度的泰勒展开去拟合Hessian的逆 B_k ，使得其带有曲率信息)

Which kind of curvature info is needed?

$$\nabla f(x) - \nabla f(y) \approx H(x - y)$$

$$\Delta x = x^{k+1} - x^k$$

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

Secant
condition



$$\Delta g \approx M^{k+1} \Delta x$$

or

$$\Delta x \approx B^{k+1} \Delta g, \quad M^{k+1} B^{k+1} = I$$

Construct B directly from gradient/variable diff

如何来增量式更新 B_k ?

Infinitely many B satisfy the secant condition?

$$\Delta x = B^{k+1} \Delta g$$

How to choose the best one?

- Assume the function is convex
- Assume an initial guess is given

$$\begin{array}{ccc} \min_B \|B - B^k\|^2 & \xrightarrow{\text{scale invariant}} & \min_B \|H^{\frac{1}{2}}(B - B^k)H^{\frac{1}{2}}\|^2 \\ \text{s.t. } B = B^T & & \text{s.t. } B = B^T \\ \Delta x = B\Delta g & & \Delta x = B\Delta g \end{array}$$

$$H = \int_0^1 \nabla^2 f[(1-\tau)x^k + \tau x^{k+1}] d\tau$$

1. 代价函数为了使得能够利用所有的 $\Delta x, \Delta g$ 对, 并且使得 B 变动最小
2. 1. 使得 B 变动最小含义就是每次在利用当前点的曲率信息的时候尽可能的利用之前的点的曲率信息

- 严格凸函数更新自然满足 B 正定性

Recover the curvature (approx. inv. Hessian) from gradients

Broyden-Fletcher-Goldfarb-Shanno (BFGS) update

$$B^{k+1} = \left(I - \frac{\Delta x \Delta g^T}{\Delta g^T \Delta x} \right) B^k \left(I - \frac{\Delta g \Delta x^T}{\Delta g^T \Delta x} \right) + \frac{\Delta x \Delta x^T}{\Delta g^T \Delta x}$$

$$B^0 = I, \Delta x = x^{k+1} - x^k, \Delta g = \nabla f(x^{k+1}) - \nabla f(x^k)$$

What conditions should $M = 1/B$ satisfy?

- It should not need full 2nd-order derivatives ✓
- The linear equations have closed-form solutions ✓
- It should be lightweight and compact to store ✓
- It must preserve descent directions ? strict convexity (necessary?)
- It should contain curvature information ✓

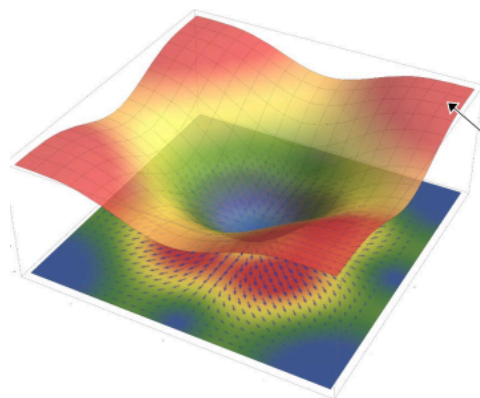
BFGS update preserves PD if $\Delta g^T \Delta x > 0$

伪代码:

BFGS method for **strictly convex** functions

```
initialize  $x^0, g^0 \leftarrow \nabla f(x^0), B^0 \leftarrow I, k \leftarrow 0$ 
while  $\|g^k\| > \delta$  do
   $d \leftarrow -B^k g^k$ 
   $t \leftarrow$  backtracking line search (Armijo)
   $x^{k+1} \leftarrow x^k + td$ 
   $g^{k+1} \leftarrow \nabla f(x^{k+1})$ 
   $B^{k+1} \leftarrow \text{BFGS}(B^k, g^{k+1} - g^k, x^{k+1} - x^k)$ 
   $k \leftarrow k + 1$ 
end while
return
```

缺点：没有梯度单调性/非凸/非光滑



Restrictions:

- Strict gradient monotonicity does not hold in general
- Curvature info far from optimum is less important
- Iteration cost is quadratic in dim
- Applicability to nonconvex function remains to be verified
- Applicability to nonsmooth function remains to be verified

- negative curvature
- nonconvex
- non-monotonic grad

How to apply to more **nonconvex/nonsmooth** functions?

3. Quasi-Newton Methods (非凸光滑函数)

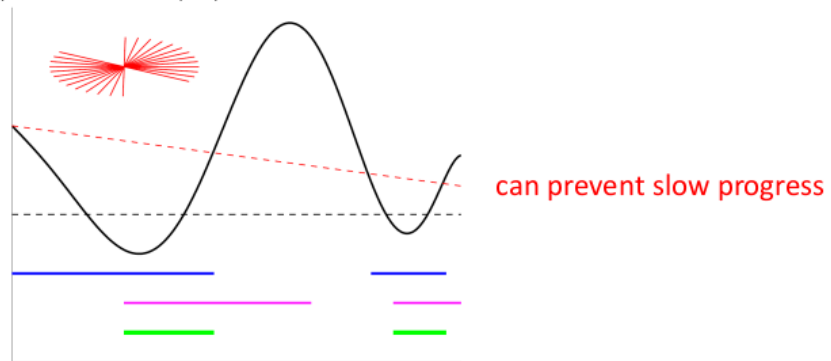
- Weak Wolfe condition (常用)

$\Delta g^T \Delta x > 0$? Ensure PD of approximation: Wolfe conditions!

1. weak Wolfe conditions $0 < c_1 < c_2 < 1$ typically $c_1 = 10^{-4}$, $c_2 = 0.9$

$$f(x^k) - f(x^k + \alpha d) \geq -c_1 \cdot \alpha d^T \nabla f(x^k) \quad \text{sufficient decrease condition}$$

$$d^T \nabla f(x^k + \alpha d) \geq c_2 \cdot d^T \nabla f(x^k) \quad \text{curvature condition}$$



保证函数下降充分的快，下降的progress更大(文献证明 Δg 和 Δx 内积在满足wolfe condition时是大于0的，从而保证B是正定的，从而保证d是函数descent方向)

1. 充分下降条件 (Armijo condition)
2. 曲率条件 (curvature condition)

• Strong Wolfe condition

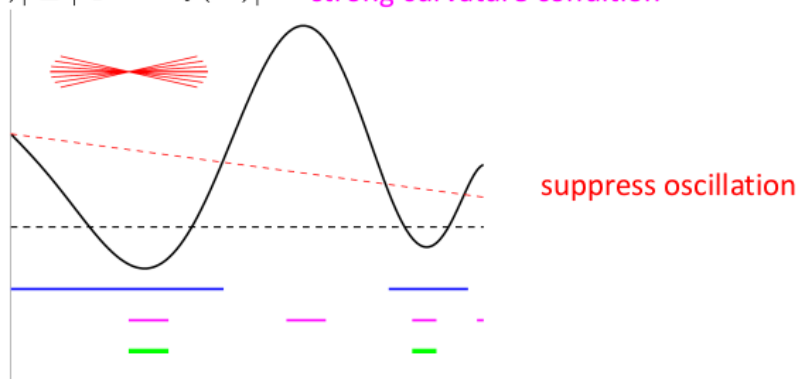
直观理解：希望下次的搜索到山谷里，把斜率压得更平（抑制震荡效果）

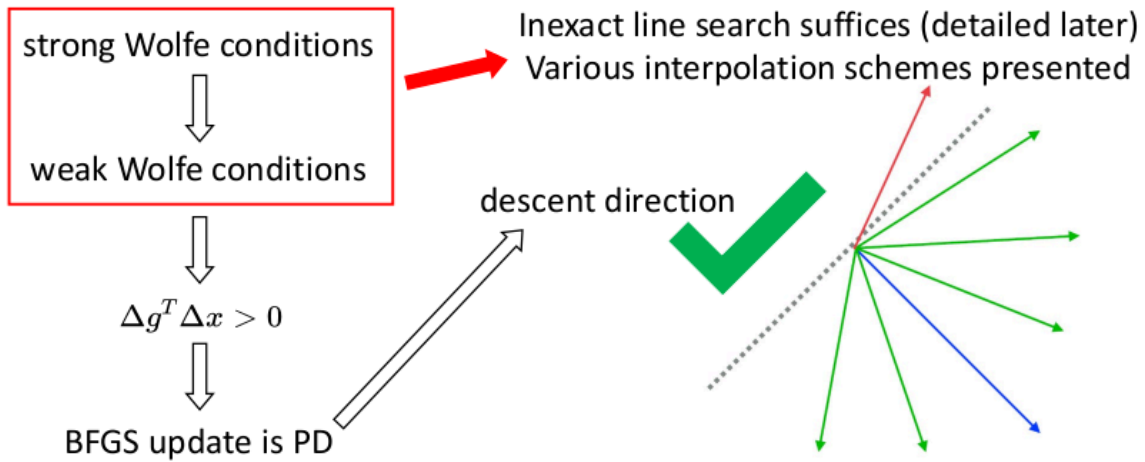
$\Delta g^T \Delta x > 0$? Ensure PD of approximation: Wolfe conditions!

2. strong Wolfe conditions $0 < c_1 < c_2 < 1$ typically $c_1 = 10^{-4}$, $c_2 = 0.9$

$$f(x^k) - f(x^k + \alpha d) \geq -c_1 \cdot \alpha d^T \nabla f(x^k) \quad \text{sufficient decrease condition}$$

$$|d^T \nabla f(x^k + \alpha d)| \leq c_2 \cdot |d^T \nabla f(x^k)| \quad \text{strong curvature condition}$$





- cautious-BFGS

Wolfe conditions **cannot guarantee** the convergence of BFGS

The **cautious update** (Li and Fukushima 2001) with mild conditions is needed.

$$B^{k+1} = \begin{cases} \left(I - \frac{\Delta x \Delta g^T}{\Delta g^T \Delta x} \right) B^k \left(I - \frac{\Delta g \Delta x^T}{\Delta g^T \Delta x} \right) + \frac{\Delta x \Delta x^T}{\Delta g^T \Delta x} & \text{if } \Delta g^T \Delta x > \epsilon \|g_k\| \Delta x^T \Delta x, \epsilon = 10^{-6} \\ B^k & \text{otherwise} \end{cases}$$

The cautious update guarantees zero grad to be an accumulation point if

- The func has **bounded sub-level sets**
- The func has **Lipschitz continuous grad**

伪代码:

BFGS method for **possibly nonconvex** functions

```

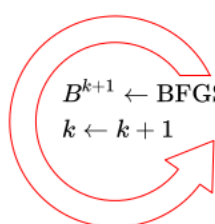
initialize  $x^0, g^0 \leftarrow \nabla f(x^0), B^0 \leftarrow I, k \leftarrow 0$ 
while  $\|g^k\| > \delta$  do
     $d \leftarrow -B^k g^k$ 
     $t \leftarrow$  inexact line search (Wolfe)
     $x^{k+1} \leftarrow x^k + td$ 
     $g^{k+1} \leftarrow \nabla f(x^{k+1})$ 
     $B^{k+1} \leftarrow$  Cautious-BFGS( $B^k, g^{k+1} - g^k, x^{k+1} - x^k$ )
     $k \leftarrow k + 1$ 
end while
return
```

4. L-BFGS算法（非凸光滑函数）

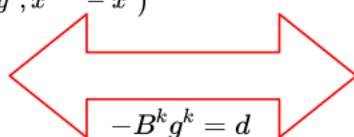
Limited-memory BFGS (**L-BFGS**): do not store B^k explicitly

$$s^k = x^{k+1} - x^k, \quad y^k = g^{k+1} - g^k, \quad \rho^k = 1 / \langle s^k, y^k \rangle \quad \text{where} \quad \langle a, b \rangle := a^T b$$

- Instead we store up to m (e.g., $m = 30$) values of s^k, y^k, ρ^k



repeat for m times $O(mn^2)$

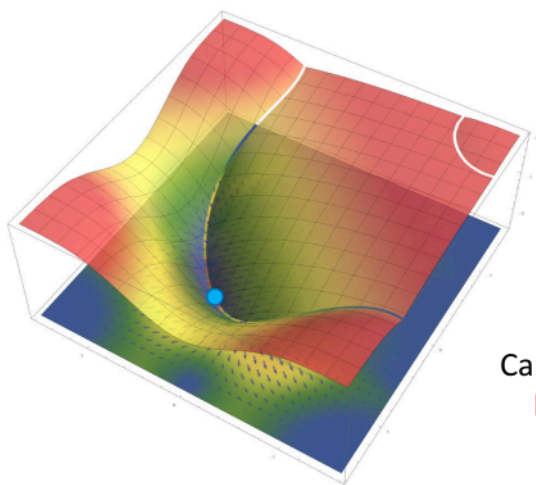


```

 $d \leftarrow g^k$   $O(mn)$ 
for  $i = k - 1, k - 2, \dots, k - m$ 
     $\alpha^i \leftarrow \rho^i \langle s^i, d \rangle$ 
     $d \leftarrow d - \alpha^i y^i$ 
end (for)
 $\gamma \leftarrow \rho^{k-1} \langle y^{k-1}, y^{k-1} \rangle$ 
 $d \leftarrow d / \gamma$ 
for  $i = k - m, k - m + 1, \dots, k - 1$ 
     $\beta \leftarrow \rho^i \langle y^i, d \rangle$ 
     $d \leftarrow d + s^i (\alpha^i - \beta)$ 
end (for)
return search direction  $d$ 
    
```

- 维护一个时间滑动窗口
- 滑动窗口入队时可以检验是否满足梯度收敛于0的条件（保证descent方向和收敛性）

5. L-BFGS算法（非凸非光滑函数）



What if functions are nonsmooth?

Troubles with nonsmoothness :

- Gradient may not exist 梯度不存在
- Negative sub-grad does not descent 次梯度not descent
- Curvature can be very large 曲率很大

Can we apply **L-BFGS** to nonsmooth function?

Practically yes! (Lewis and Overton 2013)

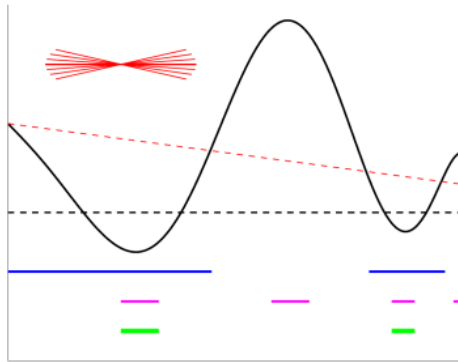
- strong wolfe conditions（不能用）

Key issues when applying L-BFGS to nonsmooth functions

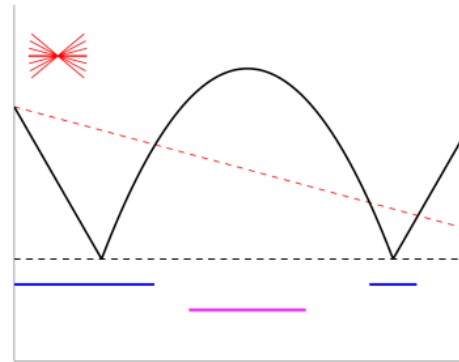
Recall: **strong Wolfe conditions**

$$f(x^k) - f(x^k + \alpha d) \geq -c_1 \cdot \alpha d^T \nabla f(x^k) \quad \text{sufficient decrease condition}$$

$$|d^T \nabla f(x^k + \alpha d)| \leq |c_2 \cdot d^T \nabla f(x^k)| \quad \text{strong curvature condition}$$



(a) Smooth $\phi(\alpha)$.



(b) nonsmooth $\phi(\alpha)$. **conditions fail!**

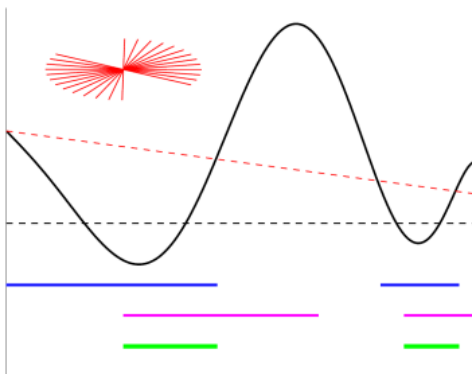
- weak wolfe conditions(适用!!)

Key issues when applying L-BFGS to nonsmooth functions

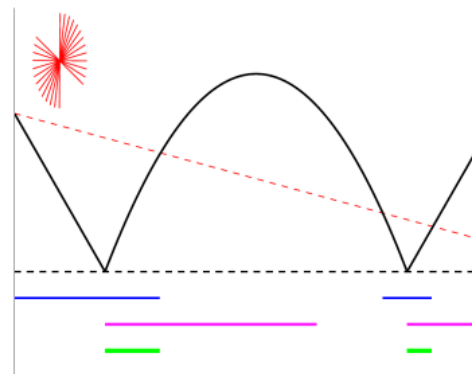
Recall: **weak Wolfe conditions**

$$f(x^k) - f(x^k + \alpha d) \geq -c_1 \cdot \alpha d^T \nabla f(x^k) \quad \text{sufficient decrease condition}$$

$$d^T \nabla f(x^k + \alpha d) \geq c_2 \cdot d^T \nabla f(x^k) \quad \text{curvature condition}$$



(a) Smooth $\phi(\alpha)$.



(b) nonsmooth $\phi(\alpha)$. **conditions work!**

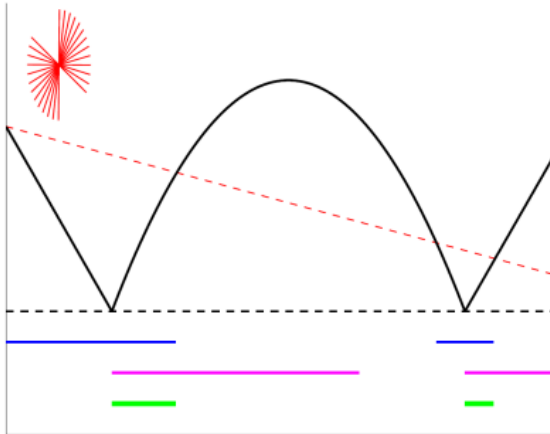
- Lewis & Overton Line search

Quasi-Newton Methods

weak Wolfe conditions should be used
for nonsmooth functions

$$S(\alpha) : f(x^k) - f(x^k + \alpha d) \geq -c_1 \cdot \alpha d^T \nabla f(x^k)$$

$$C(\alpha) : d^T \nabla f(x^k + \alpha d) \geq c_2 \cdot d^T \nabla f(x^k)$$



Lewis & Overton line search:

- weak Wolfe conditions
- no interpolation used

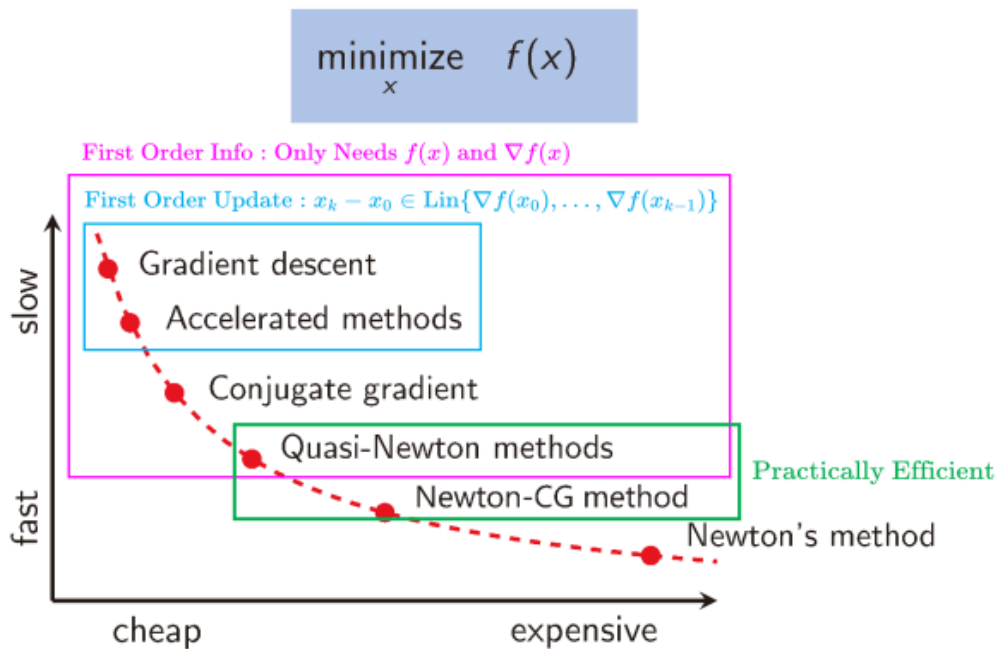
```

l ← 0
u ← +∞
α ← 1
repeat
  if S(α) fails
    u ← α
  else if C(α) fails
    l ← α
  else
    return α
  if u < +∞
    α ← (l + u)/2
  else
    α ← 2l
end (repeat)
    
```

6. Newton-CG方法(LCG)

- other method to avoid hessian calculate

Is quasi-Newton class the only way to avoid Hessian?



Complexity of function/gradient/Hessian

$$f: \mathbb{R}^n \rightarrow \mathbb{R} \quad = \quad \nabla f = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{pmatrix} \quad < \quad \nabla^2 f = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{pmatrix}$$

Function Gradient Hessian

About computational complexity :

- Computing gradient is at least $O(n)$
- Computing gradient is **as easy as** computing the function
- Computing Hessian is at least $O(n^2)$ whose inversion is about $O(n^3)$

不计算Hessian逆, 通过 $O(N)$ 计算Hessian向量

推导: $h(\alpha) = \nabla f(x + \alpha\xi)$ 在 $\alpha = 0$ 处泰勒展开

Complexity of function/gradient/Hessian-vec product

$$f: \mathbb{R}^n \rightarrow \mathbb{R} \quad = \quad \nabla f = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{pmatrix} \quad = \quad (\nabla^2 f)\xi = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{pmatrix} \begin{pmatrix} \xi_1 \\ \xi_2 \\ \vdots \\ \xi_n \end{pmatrix}$$

Function Gradient Hessian-vec product

Even if exact value is not available, just compute $(\nabla^2 f)\xi \approx \frac{\nabla f(x + \delta\xi) - \nabla f(x)}{\delta}$

However, we need **inverse-Hessian-vec** product.

What can be exploited via Hessian-vec product?

关键: 已知 $\gamma(x) = Ax$, 如何求 x

The most prominent large-scale linear solver: (Linear) Conjugate Gradient Method

$$Ax = b$$

$A \in \mathbb{R}^{n \times n}$: a known PD mat
 $x \in \mathbb{R}^n$: an unknown vec
 $b \in \mathbb{R}^n$: a known vec

Solve $Ax = b$ via the interface $\gamma(x) := Ax$

转换：二次函数求极小值的问题

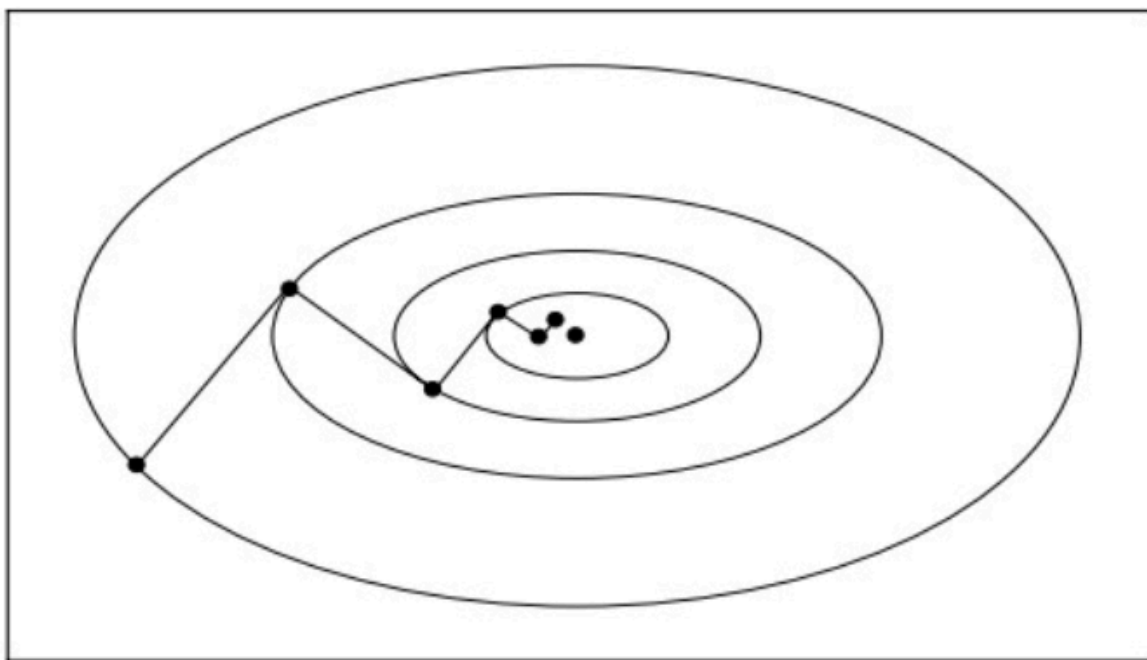
(Linear) Conjugate Gradient Method

$$\begin{array}{ccc} Ax = b & \longleftrightarrow & \arg \min_x f(x) = \frac{1}{2} x^T A x - b^T x \\ \uparrow & & \nabla f(x) = Ax - b = 0 \end{array}$$

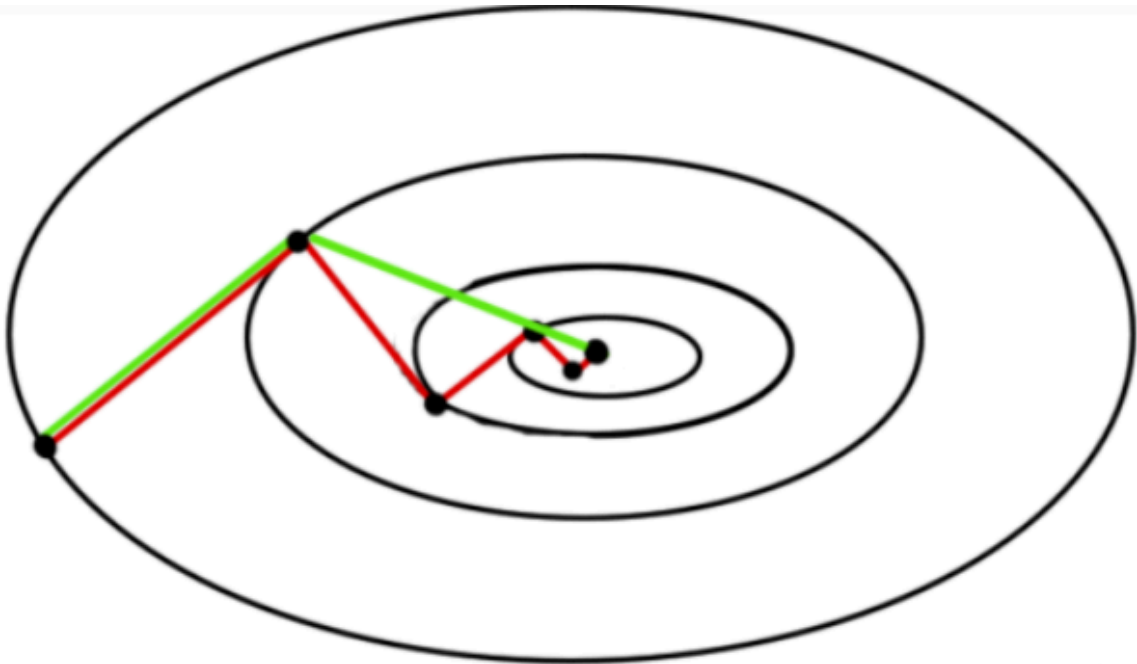
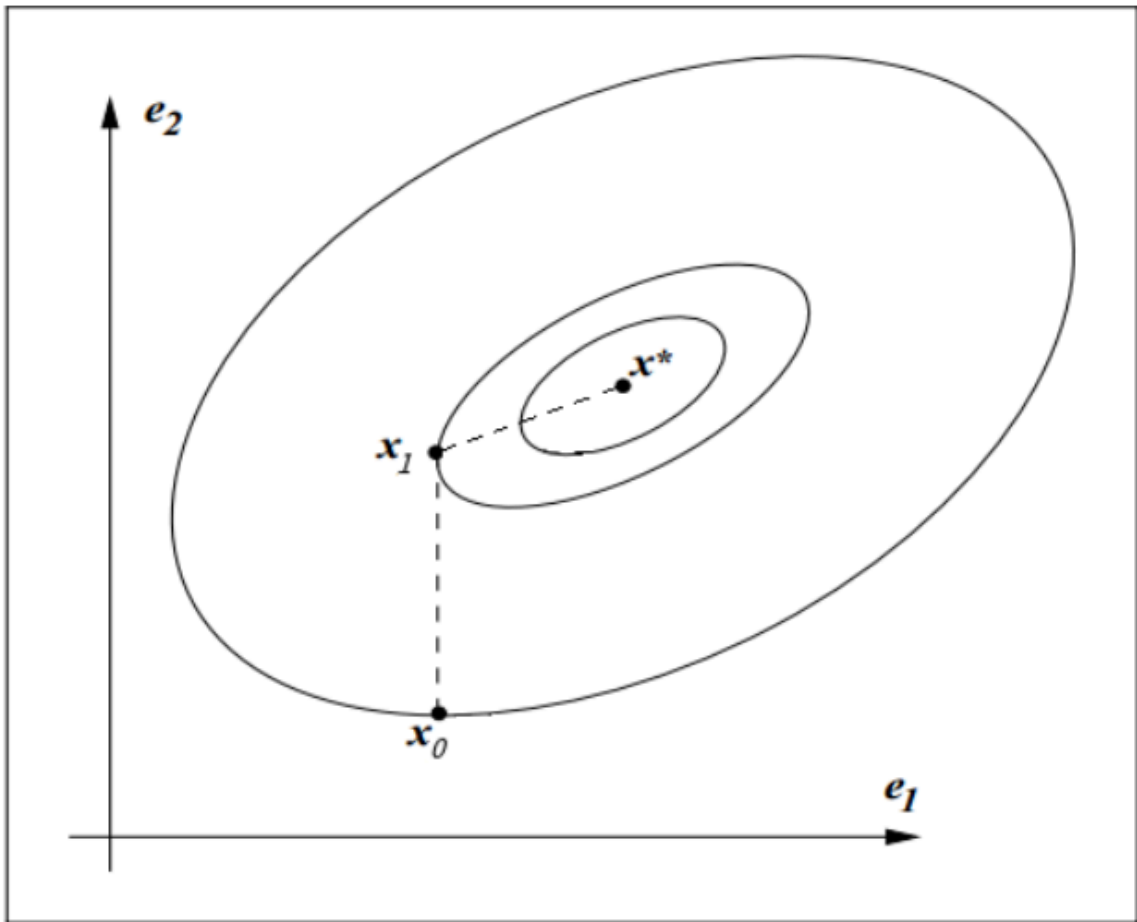
考虑使用：steepest gradient descent & newton's method, 牛顿方法不可避免求 Hessian逆, 所以考虑最速梯度下降法

问题：如何使最速梯度下降快速的迭代？

exact line search（可能会来回呈之字型震荡）



期望效果：二维函数，两步收敛（关键是第二步迭代的方向直接指向minima，由此推两个方向的关系），可以推出两个向量 x_0 和 x_1 关于 A 正交，前提是 x_0 处采用exact line search, x_1 是 x_0 exact line search搜索到的结果，推广到N维空间



关键：在 n 维空间找 n 个关于 A 正交的向量作为迭代方向

(Linear) Conjugate Gradient Method

If the linearly independent vector is chosen as

$$\mathbf{v}^k = \mathbf{b} - \mathbf{A}\mathbf{x}^k$$

The history-dependent coefficients in conjugate G-S process becomes

$$\text{proj}_{\mathbf{u}^j}(\mathbf{v}^k) = 0, \forall j \leq k-2$$

Thus we get a short recurrence

$$\mathbf{u}^k = \mathbf{v}^k - \sum_{j=1}^{k-1} \text{proj}_{\mathbf{u}^j}(\mathbf{v}^k) = \mathbf{v}^k + \frac{\|\mathbf{v}^k\|^2}{\|\mathbf{v}^{k-1}\|^2} \mathbf{u}^{k-1}$$

Given n mutually conjugate directions

$$\mathbf{u}^1, \mathbf{u}^2, \dots, \mathbf{u}^n$$

Successively conduct closed-form exact line search

$$\alpha \leftarrow \frac{\mathbf{b}^T \mathbf{u}^k - (\mathbf{x}^k)^T \mathbf{A} \mathbf{u}^k}{(\mathbf{u}^k)^T \mathbf{A} \mathbf{u}^k} \rightarrow \text{interface } \gamma(\mathbf{u}^k)$$

$$\mathbf{x}^{k+1} \leftarrow \mathbf{x}^k + \alpha \mathbf{u}^k$$

Obtain an exact solution in n steps

迭代方法（CG）解线性方程组伪代码：

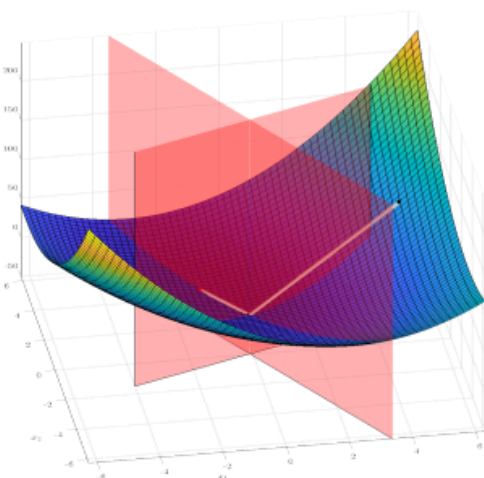
(Linear) Conjugate Gradient Method

```

 $x^1 \leftarrow x_{ini}$ 
 $\mathbf{v}^1 \leftarrow \mathbf{b} - \mathbf{A}\mathbf{x}^1$ 
 $\mathbf{u}^1 \leftarrow \mathbf{v}^1$ 
 $k \leftarrow 1$ 
while  $\mathbf{v}^k \neq 0$ 
   $\alpha \leftarrow \|\mathbf{v}^k\|^2 / (\mathbf{u}^k)^T \mathbf{A} \mathbf{u}^k$ 
   $\mathbf{x}^{k+1} \leftarrow \mathbf{x}^k + \alpha \mathbf{u}^k$ 
   $\mathbf{v}^{k+1} \leftarrow \mathbf{v}^k - \alpha \mathbf{A} \mathbf{u}^k$ 
   $\beta \leftarrow \|\mathbf{v}^{k+1}\|^2 / \|\mathbf{v}^k\|^2$ 
   $\mathbf{u}^{k+1} \leftarrow \mathbf{v}^{k+1} + \beta \mathbf{u}^k$ 
   $k \leftarrow k + 1$ 
end (while)
return  $\mathbf{x}^k$ 
  
```

float-point precision:
 $\|\mathbf{v}^k\| \geq \epsilon$

interface $\gamma(\mathbf{u}^k)$

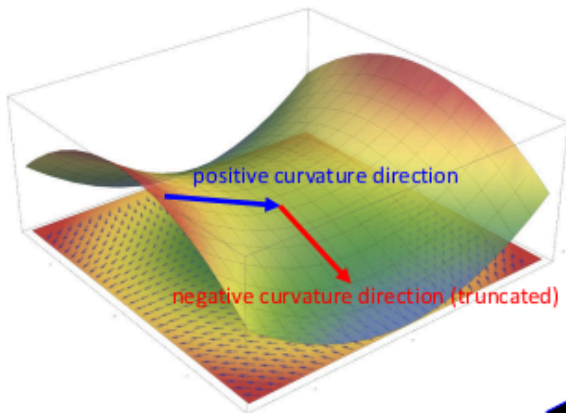


If exact arithmetic is used, CG always finds an exact solution in n steps

牛顿-CG方法用来每次更新 \mathbf{x} 时确定最佳的更新方向（求解 $\mathbf{H}\mathbf{d} = -\mathbf{g}$ 方程组）：

Newton-CG Method

Newton-CG Method



Truncated CG:

- Start with zero as the initial direction
- Truncating negative curvature direction

```

 $x^1 \leftarrow x_{ini}, k \leftarrow 1$ 
while  $\|\nabla f(x^k)\| > \delta$ 
     $\epsilon^k \leftarrow \min(1, \|\nabla f(x^k)\|/10)$ 
     $d^1 \leftarrow 0, v^1 \leftarrow -\nabla f(x^k)$ 
     $u^1 \leftarrow v^1, j \leftarrow 1$ 
    while  $\|v^j\| > \epsilon^k \|\nabla f(x^k)\|$ 
        if  $(u^j)^T \nabla^2 f(x^k) u^j \leq 0$ 
            if  $j = 1$ 
                 $d^j \leftarrow -\nabla f(x^k)$ 
            end (if)
            break
        end (if)
         $\alpha \leftarrow \|v^j\|^2 / (u^j)^T \nabla^2 f(x^k) u^j$ 
         $d^{j+1} \leftarrow d^j + \alpha u^j$ 
         $v^{j+1} \leftarrow v^j - \alpha \nabla^2 f(x^k) u^j$ 
         $\beta \leftarrow \|v^{j+1}\|^2 / \|v^j\|^2$ 
         $u^{j+1} \leftarrow v^{j+1} + \beta u^j$ 
         $j \leftarrow j + 1$ 
    end (while)
     $\alpha \leftarrow \text{backtracking line search (Armijo)}$ 
     $x^{k+1} \leftarrow x^k + \alpha d^j$ 
end (while)
    
```

1. 注意近似算 $\nabla^2 f(x_k) u^j$ 很容易出现数值不稳定情况，所以这里只在每次CG迭代外部计算一次 $\nabla^2 f(x_k)$ ，CG迭代内部只需要进行矩阵乘法
2. 整体过程参考最速梯度下降和牛顿方法，区别是不求Hessian逆而通过迭代的方法求解线性方程组，实测效率提升很大，甚至比BFGS快