

Unconstrained Minimization

(1) unconstrained minimization

minimize $f(x)$

↳ convex and differentiable

$\Rightarrow x^*$ is optimal $\Leftrightarrow \nabla f(x^*) = 0$ optimality condition

iterative algorithm $x^{(0)}, x^{(1)}, \dots \in \text{dom} f$. $f(x^{(k)}) \rightarrow p^*$ as $k \rightarrow \infty$

converge if $f(x^{(k)}) - p^* \leq \epsilon$

① initial point

要求: $\{x^{(0)} \in \text{dom} f$

sublevel set $S = \{x \in \text{dom} f \mid f(x) \leq f(x^{(0)})\}$ should be closed.

↳ satisfied for all $x^{(0)} \in \text{dom} f$ if f is closed (its sublevel sets are closed)

Definition: a function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is closed if for each α , the sublevel set

$\{x \in \text{dom} f \mid f(x) \leq \alpha\}$ is closed. then the function is closed.

properties: ① if f is a continuous function and $\text{dom} f$ is closed, then

f is closed

\mathbb{R}^n is closed.

② if f is a continuous function and $\text{dom} f$ is open, then f is

closed (iff) it converges to ∞ along every sequence converging to a boundary point of $\text{dom} f$.

② Strong convexity

f is strongly convex on S if there exists an $m \geq 0$ such that:

$\nabla^2 f(x) \succeq mI \rightarrow$ remember f is convex iff $\nabla^2 f(x) \succeq 0$

f is strictly convex iff $\nabla^2 f(x) \succ 0$

↳ 意味着 $\nabla^2 f(x)$ 的最小特征值不小于 m .

Taylor 二次展开: for $y, x \in S$

$$f(y) = f(x) + \nabla f(x)^T (y-x) + \frac{1}{2} (y-x)^T \nabla^2 f(x) (y-x)$$

$$\text{由于 } \nabla^2 f(x) \succeq mI \rightarrow \text{RHS} \geq f(x) + \nabla f(x)^T (y-x) + \frac{m}{2} (y-x)^T I (y-x)$$

$$= f(x) + \nabla f(x)^T (y-x) + \frac{m}{2} \|y-x\|_2^2 \quad \text{norm 2: } \|x\|_2 = \sqrt{x^T x}$$

联系: convex function 的 1st-order condition: $f(y) \geq f(x) + \nabla f(x)^T (y-x)$

\Rightarrow 可以说现在 $f(y) \geq f(x) + \nabla f(x)^T (y-x) + \frac{m}{2} \|y-x\|_2^2$ is a better lower bound on $f(y)$ (if $m > 0$)

28

$$f(y) \geq f(x) + \nabla f(x)^T(y-x) + \frac{m}{2} \|y-x\|_2^2$$

convex function for $y \Rightarrow y = x - \frac{\nabla f(x)}{m}$ 时有最小值

$$\begin{aligned} \Rightarrow \text{RHS} &= f(x) + \nabla f(x)^T(y-x) + \frac{m}{2} \|y-x\|_2^2 \\ &\geq f(x) + \nabla f(x)^T\left(x - \frac{\nabla f(x)}{m} - x\right) + \frac{m}{2} \left\|x - \frac{\nabla f(x)}{m} - x\right\|_2^2 \\ &= f(x) - \frac{1}{2m} \|\nabla f(x)\|_2^2 \end{aligned}$$

$$\Rightarrow f(y) \geq f(x) + \frac{1}{2m} \|\nabla f(x)\|_2^2 \quad \forall y \in S$$

$$x^*: f(x) \geq p^*$$

$$\Rightarrow p^* \geq f(x) - \frac{1}{2m} \|\nabla f(x)\|_2^2$$

$$\Rightarrow f(x) - p^* \leq \frac{1}{2m} \|\nabla f(x)\|_2^2$$

optimality condition (iteration terminates if $f(x) - p^* \leq \epsilon$)

由上可知 $\|\nabla f(x)\|_2^2$ 可提供 $f(x) - p^*$ 的 upper bound.

$$\Rightarrow \|\nabla f(x)\|_2^2 \leq (2m\epsilon)^2 \Rightarrow f(x) - p^* \leq \epsilon$$

conceptual stopping criterion because m & M are rarely known

suboptimality condition
可证当 $\nabla f(x)$ 很小时, $p(x) \rightarrow p^*$. (useful as a stopping criteria)

$$\text{还可证明得} \|x - x^*\|_2 \leq \frac{1}{m} \|\nabla f(x)\|_2$$

if f is strongly convex (即 $\nabla^2 f(x) \geq mI$) $\Rightarrow \nabla^2 f(x)$ 最大特征值也有界
即 $\nabla^2 f(x) \leq MI$

$$\text{同样: 有 } f(y) \leq f(x) + \nabla f(x)^T(y-x) + \frac{M}{2} \|y-x\|_2^2 \quad \forall x, y \in S$$

$$\Rightarrow p^* \leq f(x) + \nabla f(x)^T(y-x) + \frac{M}{2} \|y-x\|_2^2 \quad \forall x, y \in S$$

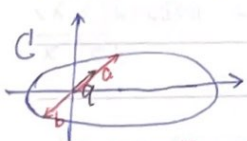
RHS achieves its minimum (for fixed x) if $y = x - \frac{\nabla f(x)}{M}$

$$\Rightarrow p^* \leq f(x) - \frac{1}{2M} \|\nabla f(x)\|_2^2$$

$$\Rightarrow MI \geq \nabla^2 f(x) \geq mI \Rightarrow \frac{1}{m} \|\nabla f(x)\|_2^2 \geq f(x) - p^* \geq \frac{1}{2M} \|\nabla f(x)\|_2^2$$

conditional # of $\nabla^2 f(x)$ (最大特征值与最小特征值之比) 上界是 $k = \frac{M}{m}$

定义 conditional #: for a convex set C , the width of C in the direction of q as $W(C, q) = \sup_{z \in C} q^T z - \inf_{z \in C} q^T z$



$$\text{minimum width: } W_{\min} = \inf_{\|q\|_2=1} W(C, q)$$

$$\text{maximum width: } W_{\max} = \sup_{\|q\|_2=1} W(C, q)$$

$$W(C, q) = q^T a - q^T b = \|a\|_2 + \|b\|_2$$

isotropic width

$$\text{Conditional \#}: \text{cond}(C) = \frac{W_{\max}^2}{W_{\min}^2}$$

\Rightarrow if conditional # of C is small \Rightarrow the set

has approximately same width in all directions (nearly spherical)

(2) Descent Method $\rightarrow f(x^{(k+1)}) < f(x^{(k)})$ except when $x^{(k)}$ is optimal

$$x^{(k+1)} = x^{(k)} + t^{(k)} \Delta x^{(k)}$$

Step size / Step length \rightarrow Step / search direction

$$f(x^{(k+1)}) \geq f(x^{(k)}) + \nabla f(x)^T (x^{(k+1)} - x^{(k)})$$

$$= f(x^{(k)}) + \nabla f(x)^T \Delta x^{(k)} + \quad (\text{convexity of } f)$$

$$\text{且 } f(x^{(k+1)}) < f(x^{(k)}) \text{ (descent method)}$$

$$\Rightarrow \nabla f(x)^T \Delta x^{(k)} < 0$$

$\Delta x^{(k)}$ is a descent direction

① Line search 找 t

1. Exact line search.

$$t = \arg \min_{s \geq 0} f(x + s \Delta x)$$

2. Backtracking Line Search (to reduce f "enough")

Start with $t=1, \alpha \in (0, \frac{1}{2}), \beta \in (0, 1)$

$$\text{while } f(x + t \Delta x) > f(x) + \alpha \nabla f(x)^T \Delta x, \quad t := \beta t$$

for fixed x , $f(x + t \Delta x)$ is a line in f with single variable t .
and fixed Δx , $f(x)$ is convex

$f(x + t \Delta x)$ is convex on t

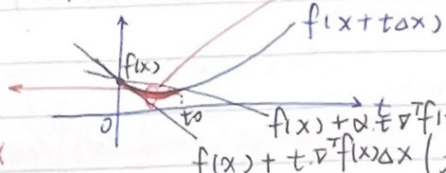
设 $t \in (0, 1)$ 内

$$f(x + t \Delta x) \leq f(x) + t \nabla f(x)^T \Delta x$$

\Rightarrow if $t \in (0, t_0)$ satisfy stopping criteria

\Rightarrow finally, $t=1$ or $t \in (\beta t_0, t_0]$

$$\text{or } t \geq \min\{1, \beta t_0\}$$



$\nabla f(x)^T \Delta x < 0$ 且 $\alpha \in (0, \frac{1}{2})$

所以 $f(x) + \alpha \nabla f(x)^T \Delta x > f(x) + t \nabla f(x)^T \Delta x$
 $t \rightarrow 0$: $f(x) \approx \text{RHS} \leq \text{LHS}$ 达到 stopping

要求 \Rightarrow Backtracking 最终必会收敛

是 $f(x + t \Delta x)$ 的一阶 Taylor 展开

$$f(x + t \Delta x) \geq f(x) + t \nabla f(x)^T \Delta x$$

② 找 Δx

1. Gradient Descent Method.

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

Convergence analysis:

① Exact Line Search

$$x^+ = x - t \nabla f(x)$$

$$f(x^+) = f(x - t \nabla f(x)) \leq f(x) - t \|\nabla f(x)\|_2^2 - \frac{Mt^2}{2} \|\nabla f(x)\|_2^2$$

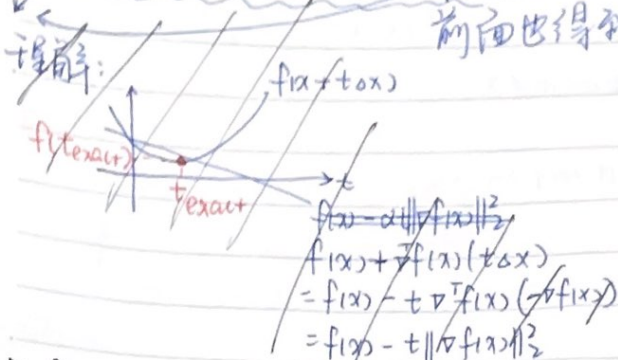
由于 line search 应找到 t 使得 $f(x^+)$ 最小

for $RHS = f(x) - t \|\nabla f(x)\|_2^2 + \frac{M}{2} t^2 \|\nabla f(x)\|_2^2$. when $t = \frac{1}{M}$

$$t = \frac{1}{M}: RHS = f(x) - \frac{1}{2M} \|\nabla f(x)\|_2^2$$

$$\Rightarrow f(t_{exact}) \leq f(x) - \frac{1}{2M} \|\nabla f(x)\|_2^2$$

前面也得到了 $f(y) \leq f(x) - \frac{1}{2M} \|\nabla f(x)\|_2^2$ 的结论



$$\Rightarrow f(x^+) \leq f(x) - \frac{1}{2M} \|\nabla f(x)\|_2^2$$

$$f(x^+) - p^* \leq f(x) - p^* - \frac{1}{2M} \|\nabla f(x)\|_2^2$$

$$\text{by } \|\nabla f(x)\|_2^2 \geq 2m(f(x) - p^*)$$

$$\text{So } f(x^+) - p^* \leq (1 - \frac{m}{M})(f(x) - p^*)$$

$$\Rightarrow f(x^{(k)}) - p^* \leq (1 - \frac{m}{M})^k (f(x^{(0)}) - p^*)$$

$$\log(\text{error of } k\text{th iteration}) \leq k \log(1 - \frac{m}{M}) + \log(\text{error of initial point})$$

lies below a line on log-linear

linear convergence

stopping criteria: $f(x^{(k)}) - p \leq \epsilon$

stop if $(1 - \frac{m}{M})^k (f(x^{(0)}) - p^*) \leq \epsilon$

$$\Rightarrow k \geq \frac{\log((f(x^{(0)}) - p^*)/\epsilon)}{\log(1/c)} \quad c = 1 - \frac{m}{M}$$

\Rightarrow we have $f(x^{(k)}) - p^* \leq \epsilon$ after at most $\frac{\log((f(x^{(0)}) - p^*)/\epsilon)}{\log(1/c)}$ iterations

for large $\frac{M}{m}$ $\frac{1}{c} = \frac{1}{1 - \frac{m}{M}} \approx \frac{M}{m}$

by $-\log x \rightarrow -x+1$ if $x \rightarrow 1$

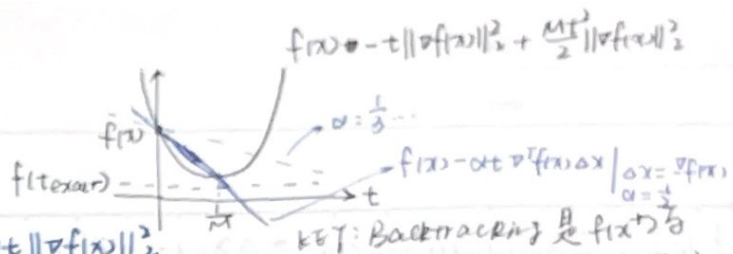
RP $-\log(1-x) \rightarrow x$ if $x \rightarrow 0$

$\Rightarrow -\log(1 - \frac{m}{M}) \rightarrow \frac{m}{M}$ if $\frac{m}{M} \rightarrow 0$

\Rightarrow if $\frac{M}{m} \rightarrow \infty$ the bound on the # of iterations increases approximately

linearly increase with $\frac{M}{m}$

同理 for Backtracking line search.



要保证 $f(t_{\text{back}}) \leq f(x) - \alpha t \|\nabla f(x)\|_2^2$
 $\therefore f(x^+) \leq f(x) - t \|\nabla f(x)\|_2^2 + \frac{M+1}{2} \|\nabla f(x)\|_2^2$
 可证 $t \in [0, \frac{1}{M}]$ 时 $\text{RHS} \leq f(x) - \alpha t \|\nabla f(x)\|_2^2$ (if $\alpha \leq \frac{1}{2}$)
 \therefore Backtracking terminates if $\begin{cases} t=1 \rightarrow f(x^+) \leq f(x) - \alpha \|\nabla f(x)\|_2^2 \\ \text{or } t \geq \frac{1}{M} \rightarrow f(x^+) \leq f(x) - \frac{\alpha}{M} \|\nabla f(x)\|_2^2 \end{cases}$ 因此要找到 α

$$\Rightarrow f(x^+) \leq f(x) - \min\{\alpha, \frac{\alpha}{M}\} \|\nabla f(x)\|_2^2$$

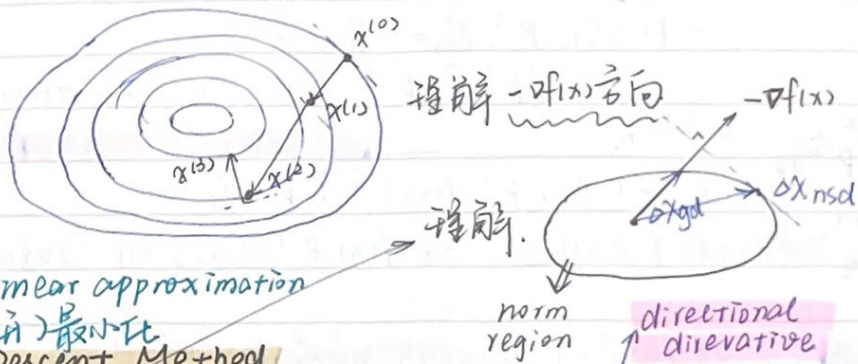
$$f(x^+) - p^* \leq f(x) - p^* - \min\{\alpha, \frac{\alpha}{M}\} \|\nabla f(x)\|_2^2$$

$$\leq (1 - \min\{2\alpha, 2\alpha \frac{M}{M}\}) (f(x^{(0)}) - p^*)$$

$$\Rightarrow f(x^{(k)}) - p^* \leq c^k (f(x^{(0)}) - p^*)$$

$$c = 1 - \min\{2\alpha, 2\alpha \frac{M}{M}\}$$

→ Gradient Descent works well if conditional # is small.



使 $f(x+v)$ 的 linear approximation (Taylor → 展并) 最小化
 2. Steepest Grad Descent Method.

linear approximation of f : $f(x+v) \approx \hat{f}(x+v) = f(x) + \nabla f(x)^T v$
 normalized steepest descent direction: $\Delta x_{\text{nsd}} = \arg \min \{ \nabla f(x)^T v \mid \|v\| \leq 1 \}$
 $\Delta x = \arg \min \{ \nabla f(x)^T v \mid \|v\| = 1 \}$

Δx_{nsd} is the direction of unit ball of $\|\cdot\|$ that extends farthest in the direction $-\nabla f(x)$

Steepest descent direction: $\Delta x_{\text{sd}} = \frac{\|\nabla f(x)\|_*}{\|\nabla f(x)\|_*^2} \Delta x_{\text{nsd}}$
 $\rightarrow \|\cdot\|_*$ is dual norm (a scalar)
 $\|a\|_* = \sup \{ x^T a \mid \|x\| = 1 \}$ $\Delta x = \sup \{ x^T a \mid \|x\| \leq 1 \}$

$\nabla f(x)^T \Delta x_{\text{sd}} = \nabla f(x)^T \cdot \frac{\|\nabla f(x)\|_*}{\|\nabla f(x)\|_*^2} \Delta x_{\text{nsd}} = -\frac{1}{\|\nabla f(x)\|_*}$
 \rightarrow 理解: $\Delta x_{\text{nsd}} = \arg \min \{ \nabla f(x)^T v \mid \|v\| = 1 \}$, $\|\nabla f(x)\|_* = \max \{ \nabla f(x)^T x \mid \|x\| = 1 \}$
 $\Rightarrow \frac{\nabla f(x)^T \Delta x_{\text{nsd}}}{\|\nabla f(x)\|_*} = \min \{ \nabla f(x)^T v \mid \|v\| = 1 \}$
 $\Rightarrow \nabla f(x)^T \|\nabla f(x)\|_* \Delta x_{\text{nsd}} = \|\nabla f(x)\|_* \cdot \min \{ \nabla f(x)^T v \mid \|v\| = 1 \} = -\|\nabla f(x)\|_* \max \{ \nabla f(x)^T v \mid \|v\| = 1 \}$
 $= -\|\nabla f(x)\|_*^2$

if $\|\cdot\|$ is Euclidean Norm $\rightarrow \Delta x_{sd} = -\nabla f(x)$
 \rightarrow Steepest Descent = Gradient Descent

2) Quadratic Norm $\|z\|_p = (z^T P z)^{1/2} = \|P^{1/2} z\|_2, P \in S_{++}^n$

$$\Delta x_{nsd} = -(\nabla f(x)^T P^{-1} \nabla f(x))^{-1/2} P^{-1} \nabla f(x)$$

$$\Delta x_{sd} = -P^{-1} \nabla f(x)$$

同理: $\Delta x_{nsd} = \arg \min \{ \nabla f(x)^T x \mid \|P^{1/2} x\|_2 = 1 \}$

let $y = P^{1/2} x$, then $x = P^{-1/2} y$

$$\Rightarrow \Delta x_{nsd} = \arg \min \{ \nabla f(x)^T P^{-1/2} y \mid \|y\|_2 = 1 \}$$

$$\Rightarrow \Delta y_{nsd} = - \frac{\nabla f(x)^T P^{-1/2}}{\|\nabla f(x)^T P^{-1/2}\|_2} = - \frac{P^{-1/2} \nabla f(x)}{\|P^{-1/2} \nabla f(x)\|_2}$$

for the Denominator: $\| \nabla f(x)^T P^{-1/2} \|_2 = \sqrt{(\nabla f(x)^T P^{-1/2})^T (\nabla f(x)^T P^{-1/2})}$

$$P \in S_{++}^n, P^{-1/2} \in S_{++}^n \Rightarrow \|P^{-1/2} \nabla f(x)\|_2 = \sqrt{\nabla f(x)^T P^{-1} \nabla f(x)}$$

$$\Rightarrow (P^{-1/2})^T = P^{-1/2} = \sqrt{\nabla f(x)^T P^{-1} \nabla f(x)}$$

$$= (\nabla f(x)^T P^{-1} \nabla f(x))^{1/2} = \|\nabla f(x)\|_{P^{-1}}$$

$$\|\nabla f(x)\|_* = \|\nabla f(x)^T P^{-1/2}\|_2$$

$$\Rightarrow \Delta y_{nsd} = -(\nabla f(x)^T P^{-1} \nabla f(x))^{-1/2} (P^{-1/2} \nabla f(x))$$

$$\Rightarrow \Delta x_{nsd} = P^{-1/2} \Delta y_{nsd} = -(\nabla f(x)^T P^{-1} \nabla f(x))^{-1/2} (P^{-1}) \nabla f(x)$$

$$\|\nabla f(x)\|_* = \max \{ \nabla f(x)^T x \mid \|x\|_P = 1 \} = \max \{ \nabla f(x)^T P^{-1/2} y \mid \|y\|_2 = 1 \}$$

$$= \|\nabla f(x)^T P^{-1/2}\|_2$$

$$\Rightarrow \Delta x_{sd} = \|\nabla f(x)\|_* \cdot \Delta x_{nsd} = \|\nabla f(x)\|_* \cdot P^{-1/2} \Delta y_{nsd}$$

$$= \|\nabla f(x)\|_* \cdot P^{-1/2} \cdot (-1) \frac{P^{-1/2} \nabla f(x)}{\|\nabla f(x)^T P^{-1/2}\|_2}$$

$$= -P^{-1} \nabla f(x)$$

\rightarrow the steepest descent method in $\|\cdot\|_P$ can be thought of as the gradient method applied to the perform after the change of coordinates $y = P^{1/2} x$

同样, Steepest Descent Method 也为 linear convergence

即: 可写成 $\|f(x^{(k)}) - p^*\| \leq C^k (\|f(x^{(0)}) - p^*\|)$

\rightarrow Steepest Descent works well ^{constant}

if the transformed problem has moderate conditional #

33

3. Newton's Method.

$$\Delta x_{nt} = -\nabla^2 f(x)^{-1} \nabla f(x)$$

→ 联系 $\Delta x_{sd} = -P^{-1} \nabla f(x) \Rightarrow$ Newton's method 也是 steepest descent

method with $\|\cdot\|$ is $\|\cdot\|_{\nabla^2 f(x)}$ 且 $P = \nabla^2 f(x)$

理解①: $\Delta x_{nt} = -\nabla^2 f(x)^{-1} \nabla f(x)$ 使 $f(x+v)$ 的 Taylor = 2 阶展开最小化

second-order Taylor approximation: $\hat{f}(x+v) = f(x) + \nabla f(x)^T v + \frac{1}{2} v^T \nabla^2 f(x) v$

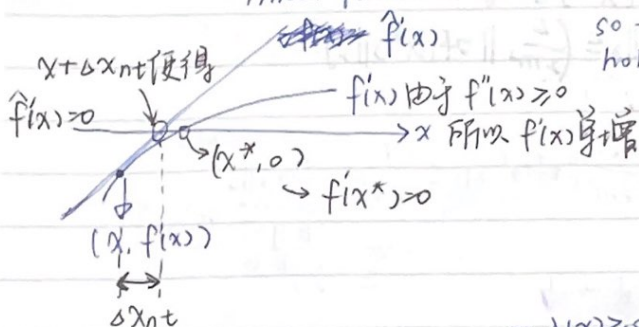
\Rightarrow if $f(x)$ is quadratic or linear then $\hat{f}(x) = f(x) \Rightarrow \Delta x_{nt} + x$ is the exact minimizer of f

理解②: 原因: we want to search a direction of v so that $\nabla f(x+v) = 0$
stopping criteria (optimality condition): $\nabla f(x^*) = 0$

$$\nabla \hat{f}(x) = \nabla f(x) + \nabla^2 f(x) v = 0 \rightarrow v = -\nabla^2 f(x)^{-1} \nabla f(x) = \Delta x_{nt}$$

→ linear function in v .

→ $\Delta x_{nt} = -\nabla^2 f(x)^{-1} \nabla f(x)$ must be added to x so that the linearized optimality condition holds.



$$\lambda(x) \geq 0, \because \nabla^2 f(x) \in S_+^n \rightarrow \nabla^2 f(x)^{-1} \in S_+^n$$

$$\Rightarrow \forall x, x^T (\nabla^2 f(x)^{-1}) x \geq 0$$

$$\text{Newton Decrement: } \lambda(x) = (\nabla f(x)^T \nabla^2 f(x)^{-1} \nabla f(x))^{\frac{1}{2}} \quad \nabla f(x)^T \nabla^2 f(x)^{-1} \nabla f(x) = 0 \text{ if } x = 0$$

$$\Rightarrow \text{if } \lambda(x) = 0$$

it means $\nabla f(x) = 0$

理解①: $\lambda(x)$ 描述 the difference between $f(x)$ and the minimum of its quadratic approximation

$$f(x) - \inf_{\Delta x} \hat{f}(x + \Delta x) = -\nabla f(x)^T \Delta x_{nt} - \frac{1}{2} \Delta x_{nt}^T \nabla^2 f(x) \Delta x_{nt}$$

$$= f(x) + \nabla f(x)^T \Delta x + \frac{1}{2} \Delta x^T \nabla^2 f(x) \Delta x$$

$$\text{令 } \Delta x_{nt} = -\nabla^2 f(x)^{-1} \nabla f(x) \text{ 代入 } \Rightarrow \text{RHS} = \frac{1}{2} \nabla f(x)^T \nabla^2 f(x)^{-1} \nabla f(x) = \frac{1}{2} \lambda(x)^2$$

$$\Rightarrow f(x) - \inf_{\Delta x} \hat{f}(x + \Delta x) = \frac{1}{2} \lambda(x)^2$$

if $\lambda(x) \downarrow \Rightarrow f(x) \rightarrow \hat{f}(x + \Delta x) \Rightarrow f(x) \rightarrow p^* \Rightarrow \lambda(x)$ can act as a stopping criteria.

$$\text{理解②: } \lambda(x) = \|\Delta x_{nt}\|_{\nabla^2 f(x)}$$

$$\|\Delta x_{nt}\|_{\nabla^2 f(x)} = (-\nabla^2 f(x)^{-1} \nabla f(x))^T \nabla^2 f(x) (-\nabla^2 f(x)^{-1} \nabla f(x))$$

$$= \nabla f(x)^T \nabla^2 f(x)^{-1} \nabla f(x) = \lambda(x)$$

$\nabla f(x)^T \Delta x_{nt}$
(directional derivative)
将 Δx_{nt} 代入, 可得

理解③: directional derivative of in the Newton's direction $\nabla f(x)^T \Delta x_{nt} = -\lambda(x)^2$

Convergence analysis

- Assume f is strongly convex $M I \leq \nabla^2 f(x) \leq m I$ for all $x \in S$
- \Rightarrow Lipschitz continuous on S 可得 $\|\nabla^2 f(x) - \nabla^2 f(y)\| \leq L \|x - y\|$

if f is a quadratic function, $\nabla^2 f(x) = \nabla^2 f(y) = a$ constant

$LHS=0 \Rightarrow L=0$ for all quadratic function

$\Rightarrow L$ measures how well f can be approximated by a quadratic model.

$\Rightarrow L$ measures the performance of Newton's method.

(Newton's method works well if L is small)

Exist constants $\eta \in (0, \frac{m^2}{L})$, $\delta > 0$ Damped Newton Phase \Rightarrow in this phase $\# \text{ of iteration} \leq \frac{f(x^{(0)}) - p^*}{\delta}$

① if $\|\nabla f(x)\| \geq \eta$, then $f(x^{k+1}) - f(x^k) \leq -\delta$

② if $\|\nabla f(x)\| < \eta$, then $\frac{L}{2m^2} \|\nabla f(x^{k+1})\| \leq (\frac{L}{2m^2} \|\nabla f(x^k)\|)^2$

Pure Newton Phase $[t=1]$
if $\|\nabla f(x^k)\| < \eta$ then $\|\nabla f(x^{k+1})\| \leq \eta$
 $\frac{L}{2m^2} \|\nabla f(x^k)\| \leq (\frac{L}{2m^2} \|\nabla f(x^0)\|)^{2^{k-1}}$
 $\leq \eta \in (0, \frac{m^2}{L})$
 $\leq (\frac{1}{2})^{2^{k-1}}$

$\Rightarrow f(x) - p^* \leq \frac{1}{2m} \|\nabla f(x)\|^2$

$\therefore f(x^k) - p^* \leq \frac{1}{2m} \|\nabla f(x^k)\|^2 \leq \frac{2m^3}{L^2} (\frac{1}{2})^{2^{k-1}}$

$\Rightarrow \log(\text{error of } k\text{th iteration}) \leq 2^{k-1} \log(\frac{1}{2}) + \log(\frac{2m^3}{L^2})$

not a linear bound versus k .

(not linear convergence)

Quadratic convergence

\rightarrow the convergence is extremely fast once the second condition is satisfied

\Rightarrow in this phase, $\# \text{ of iteration} \leq \log_2 \log_2 (\frac{E_0}{\epsilon})$

$E_0 = \frac{2m^3}{L^2}$ $f(x^{(k)}) - p^* \leq \epsilon$ (stopping criteria)

\Rightarrow overall, $\# \text{ of iterations until } f(x) - p^* \leq \epsilon$ is bounded above by $\frac{f(x^{(0)}) - p^*}{\delta} + \log_2 \log_2 (\frac{E_0}{\epsilon})$ constant

\rightarrow for different ϵ , $\log_2 \log_2 (\frac{E_0}{\epsilon})$ 值变化不大

可用 δ 估计其大小

$\Rightarrow \# \text{ of iteration} \approx \frac{f(x^{(0)}) - p^*}{\delta} + b$

13) self-concordance

前面对收敛性的分析依赖于 m, M 与 L , 而它们 are almost unknown in practice
引入 self-concordant function, 其收敛性 bound 与 m, M, L 无关且 affine invariant

Definition ①: a ^{convex} function $f: \mathbb{R} \rightarrow \mathbb{R}$ is self-concordant if $|f'''(x)| \leq 2 f''(x)^{3/2}$ for all $x \in \text{dom} f$ for function on \mathbb{R}

Standard self-concordant inequality

① 但事实上 coefficient 2 并不重要, 可认为 if a function $f(x)$ satisfies $|f'''(x)| \leq k f''(x)^{3/2}$ for $k > 0$, then 可得到 $|\tilde{f}'''(x)| \leq 2 \tilde{f}''(x)^{3/2}$ by constructing $\tilde{f}(x) = \frac{k}{4} f(x)$ 即 $\tilde{f}(x)$ is a positive-scaled version of $f(x)$

由于 positive scaling 可 preserve self-concordance (explain later) 所以 $f(x)$ is self-concordant $\Leftrightarrow \tilde{f}(x)$ self-concordant
satisfies $|\tilde{f}'''(x)| \leq 2 \tilde{f}''(x)^{3/2}$ satisfies standard self-concordant inequality

WRONG

并非所有 k 都有 \tilde{f} 为 self-concordant

系数是否为 2 有所谓, 由上分析可知 if $k \leq 2$ then $f(x)$ is self-concordant $\Rightarrow \tilde{f}(x)$ is self-concordant

② $f(x)$ is self-concordant, then $f(ax+b)$ is self-concordant

let $\tilde{f}(x) = f(ax+b)$ then \tilde{f} is self-concordant $\Leftrightarrow f$ is self-concordant

变换坐标不改变 self-concordance

Like any other positive constant k could be used, 但我不知道怎么证明

\Rightarrow self-concordance is affine invariant

Definition ③: for function on \mathbb{R}^n , a function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is self-concordant if its self-concordance holds along every line in its domain, i.e. if the function $\tilde{f}(t) = f(x + t \cdot v)$ is a self-concordant function of t all $x \in \text{dom} f$ and for all $v \in \mathbb{R}^n$

区别: $x \in \mathbb{R}, a, b \in \mathbb{R}$

Properties: (easy to +)

① scaling: if f is self-concordant and $a > 1$ then af is self-concordant

② Addition: if f_1, f_2 are self-concordant, then $f_1 + f_2$ is self-concordant

③ composition with affine function: if $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is self-concordant and $A \in \mathbb{R}^{n \times m}, b \in \mathbb{R}^n$, then $f(Ax+b)$ is self-concordant (affine invariant)

④ Composition with logarithm: let $g: \mathbb{R} \rightarrow \mathbb{R}$ be a convex function, with $\text{dom } g = \mathbb{R}_{++}$ and $|g'''(x)| \leq 3 \frac{g''(x)}{x}$ for all x . then $f(x) = -\log(-g(x)) - \log x$ is self-concordant on $\{x \mid g(x) < 0, x > 0\}$

注意: 该函数一定满足 $|g'''(x)| \leq 3 \frac{g''(x)}{x} \Rightarrow$ if $g(x)$ satisfies $|g'''(x)| \leq 3 \frac{g''(x)}{x}$

例) $g_2(x) = g_1(x) + ax^2 + bx + c$ satisfies $|g_2'''(x)| \leq 3 \frac{g_2''(x)}{x}$

例: show $f(x, y) = -\log y - \log(\log y - x)$ on $\{(x, y) \mid e^x < y\}$ is self-concordant

学习思路 ① restrict to a line

② use composition with logarithm \rightarrow 写成 $-x + \log(-\log(-))$ 形式

① restrict to a line

let $x = \hat{x} + tv$, $y = \hat{y} + tw$. where \hat{x}, \hat{y}, v, w are fixed.

$\Rightarrow f(x, y) = -\log y - \log(\log y - x) \rightarrow$ 已经是 $-\log(-) - \log(y) + \text{常数}$ 思路是写

$\Rightarrow f(\hat{x} + tv, \hat{y} + tw) = -\log(\hat{y} + tw) - \log(\log(\hat{y} + tw) - \hat{x} - tv)$ 第一个 \log 写为 $g(y)$ 且

② use composition with logarithm

(1) case I: if $w=0$

$$\text{RHS} = -\log(\hat{y}) - \log(\log \hat{y} - \hat{x} - tv)$$

\rightarrow constant

$$= -\log(at + b) + C \text{ where } a, b, C \text{ are constants}$$

because $-\log(x)$ is self-concordant, so $-\log(at + b) + C$ is also self-concordant

$\Rightarrow f(x, y)$ is self-concordant on every line

$\Rightarrow f(x, y)$ is self-concordant

(2) case II: if $w \neq 0 \Rightarrow t = \frac{y - \hat{y}}{w}$ 常用技巧

$$\Rightarrow \text{RHS} = -\log(y) - \log(\log y - x - \frac{y - \hat{y}}{w} \cdot v)$$

$$= -\log(y) - \log(\log y - a'y - b')$$

$$1 + g(y) = \log y - a'y - b'$$

$$\text{we have } g'(y) = \frac{1}{y} - a' \quad g''(y) = -\frac{1}{y^2} \quad g'''(y) = \frac{2}{y^3}$$

$$|g'''(y)| \leq 3 \frac{g''(y)}{y} \quad |g''(y)| = \frac{1}{y^2} \quad 3 \frac{g''(y)}{y} = 3 \frac{1}{y^3}$$

$$\Rightarrow g(y) \text{ satisfies } |g'''(y)| \leq 3 \frac{g''(y)}{y}$$

$\Rightarrow f(x, y) = -\log(g(y)) - \log y$ is self-concordant

37

convergence analysis

exists constants $\eta \in (0, \frac{1}{4}]$, $\delta > 0$, such that

① if $\lambda(x) > \eta$, then $f(x^{k+1}) - f(x^k) \leq -\delta$

② if $\lambda(x) \leq \eta$, then $2\lambda(x^{k+1}) \leq (2\lambda(x^k))^2$

\Rightarrow # of iteration is bounded by $\frac{f(x^{(0)}) - p^*}{\delta} + \log_2 \log_2 \left(\frac{1}{\epsilon}\right)$