

Chapter 4. Newton method.

4-1

4.1 Algorithm:

Newton method

- Initial $x = x^0$
- For $k = 0, 1, 2, \dots$
 - $\Delta x_{nt}^k = -\nabla^2 f(x^k)^{-1} \nabla f(x^k)$
 - $x^{k+1} = x^k + \eta^k \cdot \Delta x_{nt}^k$

GD: $\Delta x_{gd} = -\nabla f(x)$

Δx_{nt}^k : Newton direction.

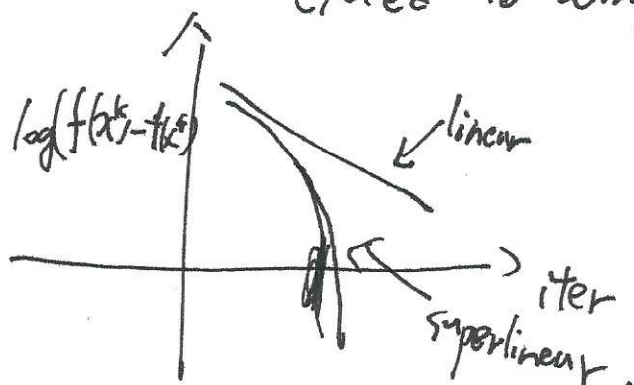
η^k : step size (usually chosen by backtracking line search + sufficient decrease condition)

Newton method: for strongly convex function:

$$\nabla^2 f(x) \geq mI, \quad m > 0$$

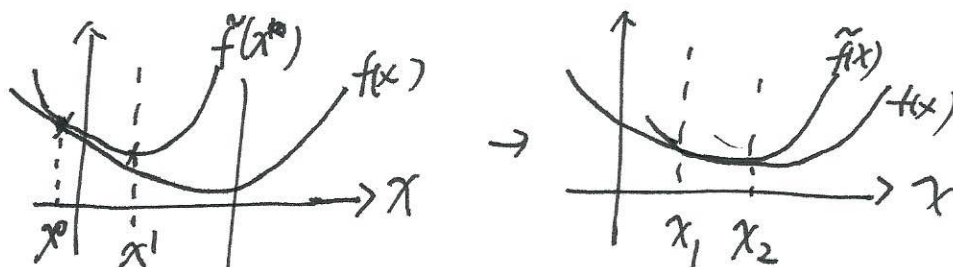
Properties: ① Faster convergence (Superlinear)

② Slower computation
(Need to compute $\nabla^2 f(x)^{-1} \nabla f(x)$)



Intuitions: Another way to explain Newton method:

4-2



At each iteration, form the "best" quadratic approximation around x^k .

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

$$\approx f(x^k + p)$$

Find the minimizer of $p^* = \operatorname{argmin} \tilde{f}(x^k + p)$

$$\nabla f(x^k) + \nabla^2 f(x^k) p^* = 0$$

$$p^* = -\nabla^2 f(x^k)^{-1} \nabla f(x^k)$$

Comparing to gradient descent:

$$\tilde{f}(x^k + p) = f(x^k) + \nabla f(x^k)^T p + \frac{1}{2} p^T I p$$

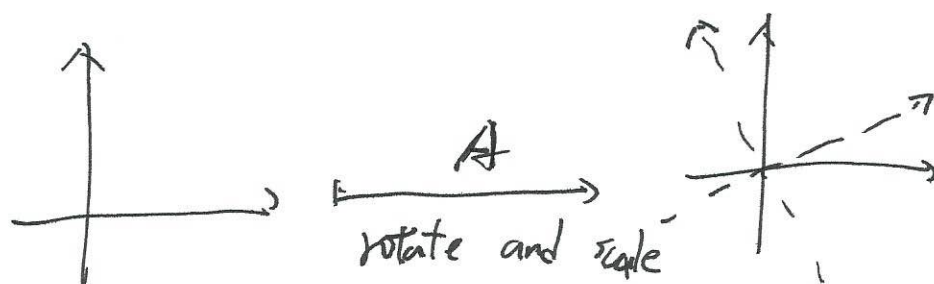
$$p^* = -\nabla f(x^k)$$

4.2 Affine invariant:

4-3

Consider two optimization problems

$$\min_x f(x) \xrightarrow{A} \min_y g(y) = f(A^{-1}y)$$



$$\begin{array}{lcl} x & \xrightarrow{\quad} & y (=Ax) \\ x^1, x^2, \dots & f(x) \xrightarrow{\quad} g(y) (=f(A^{-1}y)) & y^1, y^2, \dots \\ x^0 & \xrightarrow{\quad} & y^0 (=Ax^0) \end{array}$$

Def: The algorithm is "affine invariant" if $\forall k=1,2,\dots, y^k = Ax^k$

This means the algorithm will not change if we change the basis.

Property: Gradient Descent is "not" affine invariant, y -space

Why? Original space \xrightarrow{A} $y^1 \leftarrow y^0 - \alpha \nabla g(y^0)$

$$\begin{aligned} x^1 &\leftarrow x^0 - \alpha \nabla f(x^0) \\ &\downarrow \\ y^1 &\neq Ax^1 \end{aligned}$$

$$\begin{aligned} &= y^0 - \frac{\partial}{\partial y} f(A^{-1}y) \Big|_{y=y^0} \\ &= y^0 - A^{-T} \frac{\partial}{\partial x} f(x) \Big|_{x=x^0} \\ &= y^0 - A^{-T} \nabla_x f(x^0) \end{aligned}$$

$Ax^0 - A \nabla f(x^0)$ \neq $y^1 = Ax^0 - A^{-T} \nabla_x f(x^0)$

Thm: Newton method is affine invariant.

4-4

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

$$x^{k+1} = x^k - \eta^k \cdot \Delta x_{nt}^k$$

--- -- "I" ---

Original space

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

$$x^{k+1} = x^k - \nabla^2 f(x^k)^{-1} \cdot \nabla f(x^k)$$

y - space .

~~$$y^{k+1} = y^k - \nabla_y g(y^k)$$~~

~~$$y^{k+1} = y^k - \nabla^2 g(y^k)^{-1} \cdot \nabla g(y^k)$$~~

$$y^{k+1} = y^k - \nabla^2 g(y^k)^{-1} \cdot \nabla g(y^k)$$

$$= y^k - \left(\frac{\partial^2}{\partial y^2} f(A^{-1}y) \Big|_{y=y^k} \right)^{-1} \cdot \left(\frac{\partial}{\partial y} f(A^{-1}y) \Big|_{y=y^k} \right)$$

$$\cdot \left(\frac{\partial}{\partial y} f(A^{-1}y) \Big|_{y=y^k} \right)$$

$$= y^k - \left(A^{-T} \frac{\partial^2}{\partial x^2} f(x^k) A^{-1} \right)^{-1} \cdot \left(A^{-T} \cdot \frac{\partial}{\partial x} f(x^k) \right)$$

$$= y^k - A \cdot \nabla^2 f(x^k)^{-1} \cdot A^T \cdot \nabla f(x^k)$$

$$= y^k - A \cdot \nabla^2 f(x^k)^{-1} \cdot \nabla f(x^k)$$

$$= Ax^{k+1} - A \cdot \nabla^2 f(x^k)^{-1} \cdot \nabla f(x^k)$$

$$Ax^{k+1} = y^{k+1}$$

$(g(y) = f(A^{-1}y))$

$x = A^{-1}y$

$x^k = A^{-1}y^k$

#

4.3 Newton Decrement.

4-5

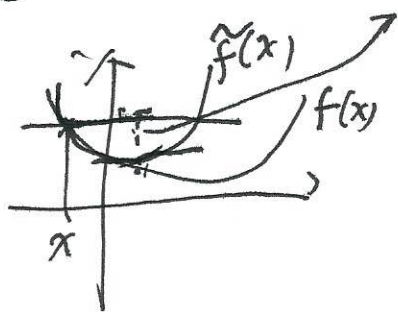
Newton Decrement:

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

① In the line search (sufficient decrease condition).

$$\begin{aligned} f(x + \eta \Delta x_{nt}) &\leq f(x) + \eta \nabla f(x)^T \Delta x_{nt} \\ &= f(x) + \underbrace{(\eta \cdot \nabla f(x)^T (\nabla^2 f(x)^{-1} \nabla f(x)))}_{\lambda(x)^2} \\ &= f(x) + c \cdot \eta \cdot \lambda(x)^2 \end{aligned}$$

② Thm: $f(x) - \inf_y \tilde{f}(y) = \frac{1}{2} \lambda(x)^2$.



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

$$\lambda + \Delta \lambda_{nt}^* = \arg \min_y \tilde{f}(y)$$

$$\begin{aligned} \tilde{f}(x + \Delta x_{nt}^*) &= f(x) - \nabla f(x)^T \nabla^2 f(x)^{-1} \nabla f(x) \\ &\quad + \frac{1}{2} \nabla f(x)^T \nabla^2 f(x)^{-1} \nabla f(x) \cdot \nabla f(x)^T \nabla f(x) \end{aligned}$$

$$\begin{aligned} &= f(x) - \nabla f(x)^T \nabla^2 f(x)^{-1} \nabla f(x) \\ &\quad + \frac{1}{2} \nabla f(x)^T \nabla^2 f(x)^{-1} \nabla f(x) \cdot \nabla f(x)^T \nabla f(x) \end{aligned}$$

$$\begin{aligned} &= f(x) - \frac{1}{2} \underbrace{\nabla f(x)^T \nabla^2 f(x)^{-1} \nabla f(x)}_{\lambda(x)^2} \\ &\quad \cdot \lambda(x)^2 \end{aligned}$$

③ $\lambda(x)$ can be used as Stopping condition.

Newton' method (with line search & stopping condition)

Initial x^0

For $k=0, 1, 2, \dots$

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

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

② Stop if $\lambda^2/2 \leq \varepsilon$ (or $\lambda^2/2 \leq \varepsilon \cdot \lambda_0^2/2$)

③ For $\eta = 1, \beta, \beta^2, \dots$

Quit if $f(x^k + \eta \Delta x_{nt}^k) \leq f(x^k) + c \cdot \eta \cdot \lambda^2$

④ Update $x^{k+1} = x^k + \eta \cdot \Delta x_{nt}^k$

GD: $0 < C < 1$

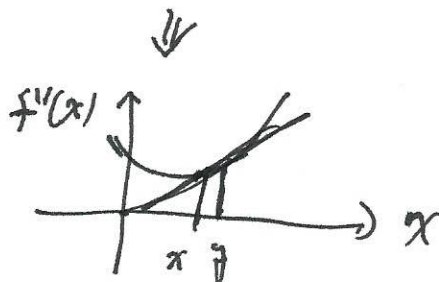
Newton: $0 < C < \frac{1}{2}$

4.4 Convergence.

Assumptions:

① f is strongly convex, twice differentiable. m -Lipchitz: $mI \leq \nabla^2 f(x) \leq M I$ ($m, M > 0$)② $\nabla^2 f(x)$ is Lipchitz continuous:

$$\|\nabla^2 f(x) - \nabla^2 f(y)\|_2 \leq L \|x - y\|_2 \quad \forall x, y.$$



~~if $\nabla^3 f(x)$~~
if $\nabla^3 f(x)$
 $\Rightarrow \|\nabla^3 f(x)\| \leq L$

Convergence properties: the iterations fall into "two" phases

- Phase I: the "global" or "damped" phase.
- Phase II: the "local" or "quadratic" phase.

There are numbers $0 < \alpha \leq m/M$ and $\gamma > 0$, st.- Phase I: if $\|\nabla f(x^k)\|_2 \geq \alpha$, then

$$f(x^{k+1}) - f(x^k) \leq -\gamma$$

- Phase II: if $\|\nabla f(x^k)\|_2 < \alpha$, then① The backtracking line search give you $\eta=1$

quadratic convergence \rightarrow ② $\frac{M}{2m^2} \|\nabla f(x^{k+1})\|_2 \leq \left(\frac{M}{2m^2} \|\nabla f(x^k)\|_2 \right)^2$

Lipchitz ~~and~~ continuous. (Clarify the definition)

Def: If a function f is Lipchitz continuous

$$\text{iff } \|f(x) - f(y)\| \leq L \cdot \|x - y\| \quad \forall x, y$$

L : Lipchitz constant, $L > 0$.

If f is differentiable:

$$f \text{ is } L\text{-Lipchitz continuous iff } \|\nabla f(x)\| \leq L \quad \forall x$$

When we say f is L -Lipchitz

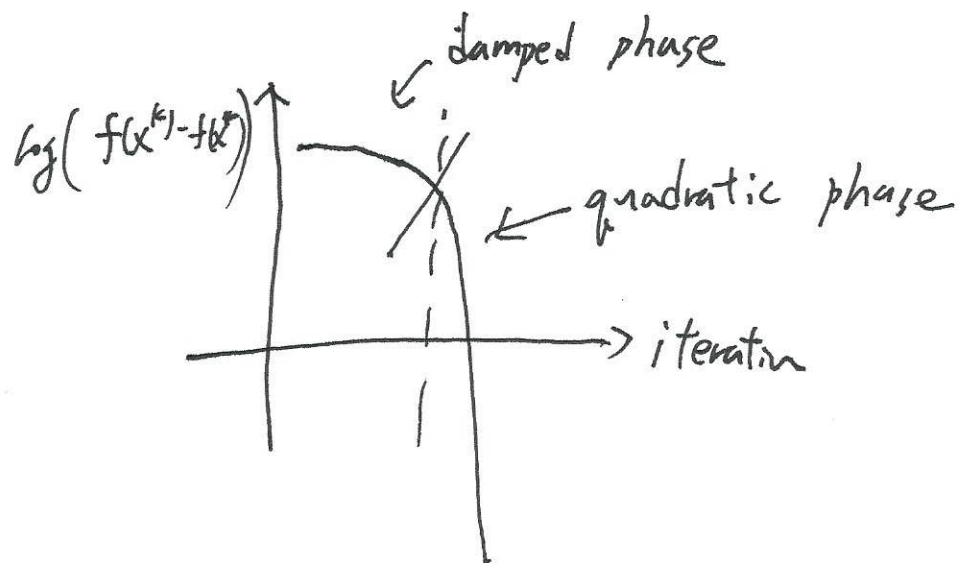
$$\Leftrightarrow \nabla f \text{ is } L\text{-Lipchitz continuous}$$

$$\checkmark \Leftrightarrow \|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\| \quad \forall x, y$$

In Newton

we want $\nabla^2 f$ is L -Lipchitz continuous

$$\checkmark \Leftrightarrow \|\nabla^2 f(x) - \nabla^2 f(y)\| \leq L \|x - y\|$$



How many iterations do we need to get
 "ε - accurate solution"?

$$\Downarrow$$

$$f(x^k) - f(x^*) \leq \varepsilon$$

Phase I: $f(x^{k+1}) - f(x^k) \leq -\gamma$

l iterations in phase I:
$$l \leq \frac{f(x^0) - f(x^*)}{\gamma}$$

Phase II: $\| \nabla f(x^k) \| < \alpha \leq \frac{m^2}{n} \quad \& \quad \frac{M}{2m^2} \| \nabla f(x^{k+1}) \| \leq \left(\frac{M}{2m^2} \| \nabla f(x^k) \| \right)^2$

\Downarrow
 a iterations

① After a iterations:

$$\frac{M}{2m^2} \| \nabla f(x^{l+a}) \| \leq \left(\frac{M}{2m^2} \| \nabla f(x^l) \| \right)^{2^a}$$

$$\leq \left(\frac{1}{2} \right)^{2^a}$$

$$\| \nabla f(x^{l+a}) \| \leq \frac{2m^2}{M} \left(\frac{1}{2} \right)^{2^a}$$

② Connect $\|\nabla f(x^{l+a})\|$ to $f(x^{l+a}) - f(x^*)$ 4-10

$$\forall y, x \quad f(y) \geq \left\{ f(x) + \nabla f(x)^T (y-x) + \frac{m}{2} \|y-x\|^2 \right\} := g(y)$$

$$\tilde{y} = \underset{y}{\operatorname{argmin}} g(y)$$

$$\nabla f(x) + m(\tilde{y} - x) = 0$$

$$\tilde{y} = x - \frac{1}{m} \nabla f(x)$$

$$\underline{f(x^*)} \geq g(x^*) \geq g(\tilde{y}) = f(x) + \nabla f(x)^T \left(-\frac{1}{m} \nabla f(x) \right) + \frac{m}{2} \left\| -\frac{1}{m} \nabla f(x) \right\|^2$$

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

$$= \underline{f(x) - \frac{1}{2m} \|\nabla f(x)\|^2}$$

$$f(x^{l+a}) - f(x^*) \leq \frac{1}{2m} \|\nabla f(x^{l+a})\|^2$$

$$\leq \frac{1}{2m} \left(\frac{2m^2}{n} \left(\frac{1}{2} \right)^{2^a} \right)^2$$

$$= \frac{2m^3}{n^2} \cdot \left(\frac{1}{2} \right)^{2^{a+1}} := \varepsilon_0 \cdot \left(\frac{1}{2} \right)^{2^a}$$

$$\text{Want: } \varepsilon_0 \cdot \left(\frac{1}{2} \right)^{2^a} < \varepsilon$$

$$\Rightarrow a \geq \log_2 \log_2 (\varepsilon_0 / \varepsilon)$$

Total number of iterations

$$l + a \leq \frac{f(x^0) - f(x^*)}{\gamma} + \underbrace{\log_2 \log_2 (\varepsilon_0 / \varepsilon)}_{\#}$$

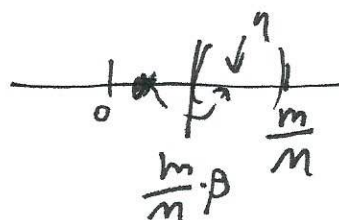
Damped Newton Phase:

Thm: $\exists \alpha \leq \frac{m^2}{m}$ such that
 if $\|\nabla f(x^k)\| \geq \alpha$, then $f(x^{k+1}) - f(x^k) \leq -\gamma$

pf: Let $x = x^k$
 $x^+ = x^{k+1}$ $x^+ = x - \eta \cdot \nabla^2 f(x)^{-1} \nabla f(x)$
 $= x + \eta \Delta x_{nr}$

$$\begin{aligned} f(x^+) &\leq f(x) + \nabla f(x)^T (\eta \Delta x_{nr}) + \frac{M}{2} \|\eta \Delta x_{nr}\|^2 \\ &= f(x) - \eta \cdot \nabla f(x)^T \nabla^2 f(x)^{-1} \nabla f(x) \\ &\quad + \frac{M}{2} \eta^2 \cdot \nabla f(x)^T \nabla^2 f(x)^{-1} \cdot \nabla^2 f(x)^{-1} \nabla f(x) \\ &\leq f(x) - \eta \cdot \lambda(x)^2 + \frac{\eta^2 M}{2m} \nabla f(x)^T \nabla^2 f(x)^{-1} \nabla f(x) \\ &= f(x) - \eta \lambda(x)^2 + \frac{M \eta^2}{2m} \lambda(x)^2 \end{aligned}$$

$$\begin{aligned} \text{When } \eta \leq \frac{m}{M} \Rightarrow f(x^+) - f(x) &\leq -\eta \lambda(x)^2 + \frac{\eta}{2} \lambda(x)^2 \\ &= -\frac{\eta}{2} \lambda(x)^2 \\ (C < \frac{1}{2}) \rightarrow &\leq -C \cdot \eta \cdot \nabla f(x)^T \Delta x_{nr} \end{aligned}$$

 step size $\geq \frac{m}{M} \beta$.

$$\begin{aligned} \Rightarrow f(x^+) - f(x) &\leq -C \cdot \eta \cdot \nabla f(x)^T \nabla^2 f(x)^{-1} \nabla f(x) \\ &\leq -C \cdot \frac{m}{M} \beta \cdot \left(\frac{1}{M} \cdot \|\nabla f(x)\|^2 \right) \\ &\leq -C \cdot \frac{m}{M^2} \beta \cdot \alpha^2 \\ &\quad \text{"} \gamma \quad \# \end{aligned}$$

Quadratic convergence phase:

4-12

① Prove $\eta=1$ after $\|\nabla f(x)\| \leq \alpha \leq m^2/n$

(see page 490-491 "Convex Optimization")

② Quadratic convergence when $\|\nabla f(x)\| \leq \alpha \leq m^2/n$

$$\text{pf: } \|\nabla f(x^+)\|_2 = \left\| \left(\nabla f(x + \alpha \Delta x_{nt}) - \nabla f(x) \right) - \underbrace{\nabla^2 f(x) \Delta x_{nt}}_{(-\nabla f(x))} \right\|_2$$

$$= \left\| \int_0^1 \nabla^2 f(x + \alpha \Delta x_{nt}) \Delta x_{nt} d\alpha - \int_0^1 \nabla^2 f(x) \Delta x_{nt} d\alpha \right\|_2$$

$$\begin{aligned} & \stackrel{\substack{L\text{-Lipsh't.} \\ \text{at } \nabla^2 f}}{\Rightarrow} = \left\| \int_0^1 \left(\nabla^2 f(x + \alpha \Delta x_{nt}) - \nabla^2 f(x) \right) \Delta x_{nt} d\alpha \right\|_2 \\ & \leq \left\| \int_0^1 L \cdot \|\alpha \cdot \Delta x_{nt}\| \cdot \|\Delta x_{nt}\| d\alpha \right\|_2 \\ & = L \|\Delta x_{nt}\|^2 \cdot \int_0^1 \alpha d\alpha \\ & = \frac{L}{2} \|\nabla^2 f(x)^T \nabla f(x)\|^2 \\ & \leq \frac{L}{2m^2} \|\nabla f(x)\|^2 \end{aligned}$$

#