

# Chapter 6 Newton method for large-scale problems.

6-1

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

Newton direction: solving linear system

$$\underbrace{\nabla^2 f(x^k)}_{n \times n} \cdot \underbrace{p^k}_{n \times 1} = \underbrace{\nabla f(x^k)}_{n \times 1}$$

"Inexact" Newton method:

Use another iterative solver for solving

$$0 = p^{k,0}, p^{k,1}, p^{k,2}, \dots, p^{k,t}, \dots$$

until  $p^{k,t}$  satisfy the stopping condition.

A typical stopping condition:

$$\|r^{k,t}\| \leq \epsilon_k \cdot \|\nabla f(x^k)\| \quad \rightarrow r^{k,0}$$

$$r^{k,t} = \nabla^2 f(x^k) \cdot p^{k,t} - \nabla f(x^k)$$

$\{\epsilon_k\}$ : forcing sequence.

Inexact Newton method:

(use CG)

For  $k=0, 1, \dots$

- solve  $\nabla^2 f(x^k) \cdot p^k = \nabla f(x^k)$  to get  $p^k$

until  $\|\nabla^2 f(x^k) p^k - \nabla f(x^k)\| \leq \epsilon_k \cdot \|\nabla f(x^k)\|$

-  $x^{k+1} = x^k + \eta^k \cdot p^k$  where  $\eta^k$  by line search.

end.

When using CG, only need  $\underbrace{\nabla^2 f(x^k)}_{n \times n} \cdot \underbrace{p^k}_{n \times 1}$  at each iteration.

Recall: convergence rate.

- Gradient descent: linear convergence.

$$\lim_{k \rightarrow \infty} \frac{f(x^{k+1}) - f(x^*)}{f(x^k) - f(x^*)} \leq C, \quad 1 > C > 0$$

- Newton method: (exact.  $\nabla^2 f(x^k)^{-1} \nabla f(x^k)$ )

$$\lim_{k \rightarrow \infty} \frac{f(x^{k+1}) - f(x^*)}{(f(x^k) - f(x^*))^2} \leq C, \quad C > 0$$

$$\lim_{k \rightarrow \infty} \frac{\|\nabla f(x^{k+1})\|}{\|\nabla f(x^k)\|^2} \leq C$$

Thm: If  $f$  is twice continuous differentiable. Then

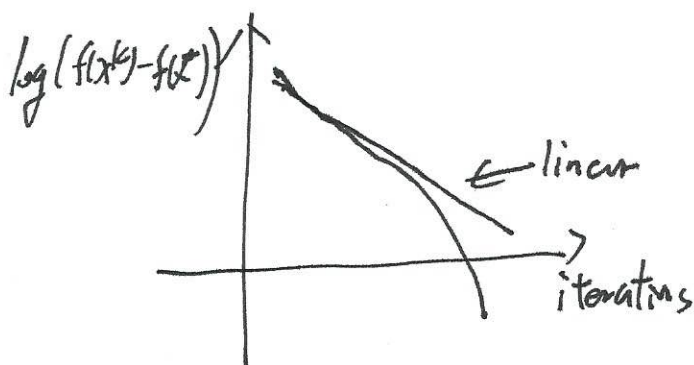
①  $\{G_k\} \rightarrow 0$ , then inexact Newton has "superlinear" convergence rate.  
e.g.  $y_k$

②  $G_k = O(\|\nabla f(x^k)\|)$ , then inexact Newton has "quadratic" convergence rate.  
 $= C \cdot \|\nabla f(x^k)\|$

Superlinear:  $\lim_{k \rightarrow \infty} \frac{f(x^{k+1}) - f(x^*)}{f(x^k) - f(x^*)} = 0$

Quadratic convergence:

$$\lim_{k \rightarrow \infty} \frac{f(x^{k+1}) - f(x^*)}{(f(x^k) - f(x^*))^2} \leq C.$$

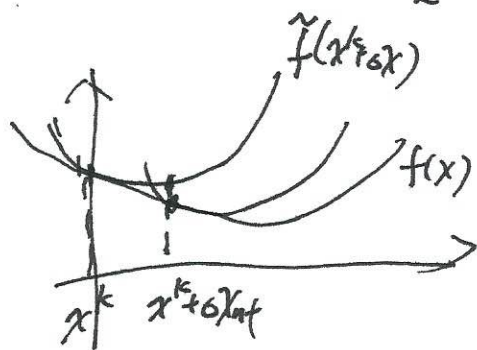


## Quasi-Newton method.

## Newton method.

- Each iteration. form a quadratic approximation

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



$$\Delta x_{nt} := \underset{\Delta x}{\operatorname{argmin}} \tilde{f}_k(x^k + \Delta x)$$

$$\nabla \tilde{f}(x^k + \Delta x_{nt}) = 0$$

$$\nabla f(x^k) + \nabla^2 f(x^k) \Delta x_{nt} = 0$$

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

Quasi-Newton: Replace  $\nabla^2 f(x^k)$  by some  $B_k$   
st  $B_k \succ 0I$

$$\Delta x'_{nt} = -B_k^{-1} \cdot \nabla f(x^k)$$

$$x^{k+1} = x^k + \eta^k \cdot \Delta x'_{nt}$$

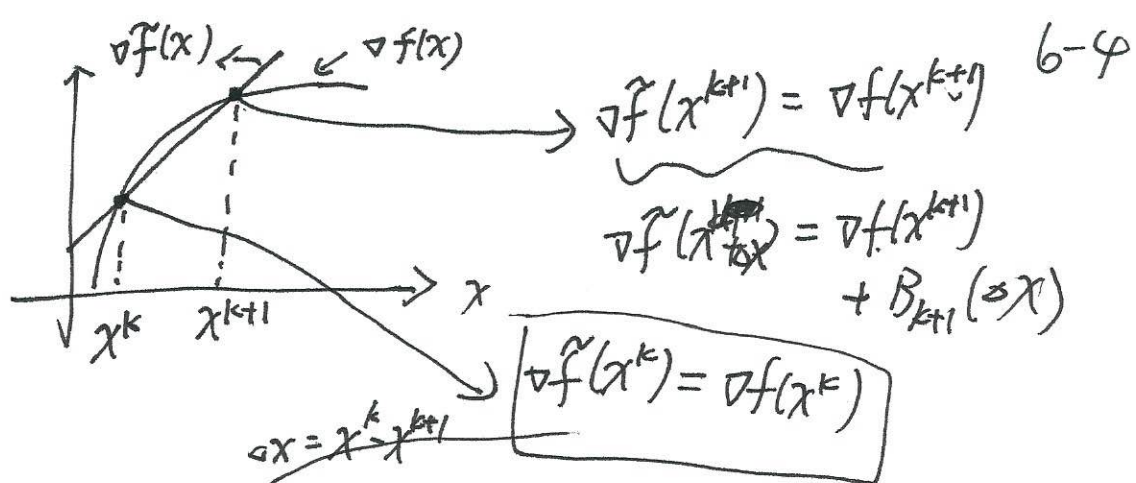
How to choose  $B_k$ ?

**BFGS**: only use the information  $\{ \nabla f(x^0), \nabla f(x^1), \nabla f(x^2) \dots; \}$   
 $\nabla f(x^k)$

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

$$\text{Want } \nabla \tilde{f}(x^k) = \nabla f(x^k)$$

the gradient for  
the approximate function  
at  $x^k$  real gradient at  $x^k$



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

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

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

$$B_{k+1} \cdot s_k = y_k \quad \text{Secant equation}$$

① There exists a (pd) solution  $B_{k+1}$  iff  $s_k^T y_k > 0$

$$(\Rightarrow) \quad s_k^T B_{k+1} s_k = s_k^T y_k > 0$$

$$(\Leftarrow) \quad B_{k+1}: n^2/2 \text{ freedom.}$$

secant equation:  $n$  equalities.

pd: constraint.

② If  $f(x)$  is strongly convex, then  $s_k^T y_k > 0$  (strictly)



BFGS:

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

$$H_{k+1} = B_{k+1}^{-1} \quad \text{Secant equation.}$$

$$B_{k+1} s_k = y_k$$

$$\boxed{H_{k+1} y_k = s_k} \quad \text{Choosing } H_{k+1} \text{ is equivalent to choosing } B_{k+1}$$

BFGS: select  $H_{k+1}$  by

$$\arg \min_{H} \|H - H_k\|_W$$

$$\text{s.t. } s_k = H y_k \quad \text{and } H > 0$$

$$\text{Solution: } H_{k+1} = (I - \rho_k s_k y_k^T) H_k (I - \rho_k y_k s_k^T) + \rho_k s_k s_k^T$$

$$\text{where } \rho_k = 1 / y_k^T s_k \quad (s_k = x^{k+1} - x^k, y_k = \nabla f(x^{k+1}) - \nabla f(x^k))$$

Usually  $H_0 = \lambda I$  for some  $\lambda > 0$ 

$$H_0 \rightarrow H_1 \rightarrow H_2 \rightarrow \dots$$

Algorithm: (BFGS)

- Initial  $x^0, H^0$ - For  $k=0, 1, \dots$ 

$$\text{Compute } p_k = - (H_k \cdot \nabla f(x^k))$$

$$\text{Set } x^{k+1} = x^k + \eta^k \cdot p_k \quad (\text{by } \eta^k \text{ by line search})$$

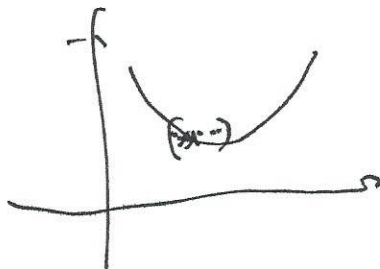
$$\text{Compute } s_k = x^{k+1} - x^k, y_k = \nabla f(x^{k+1}) - \nabla f(x^k)$$

$$\text{Compute } H_{k+1} = (I - \underbrace{\rho_k s_k y_k^T}_{V_k}) \underbrace{H_k}_{V_k} (I - \underbrace{\rho_k y_k s_k^T}_{V_k^T}) + \rho_k s_k s_k^T$$

end.

Thm = If  $f$  is twice continuously differentiable, and 6-6  
 $mI \leq \nabla^2 f(x) \leq MI$  for some  $m, M > 0$ ,  $\forall x$ .

Then  $\{x^k\} \rightarrow x^*$  superlinearly.



Computation & storage

Computation:  $\nabla f(x^k)$

$$H_{k+1} = V_k^T H_k V_k + \underbrace{\rho s_k s_k^T}_{\text{rank-1 update}}$$

$$H_{k+1} \cdot \nabla f(x^k) \leftarrow O(n^2)$$

Storage: store  $s_k, y_k \quad \forall k$  in memory.

$O(k \cdot n)$  memory

$\uparrow$   
 # iterations.

# Limited memory BFGS (L-BFGS)

6-7

Revisit the BFGS update for  $H_{k+1}$

$$H_{k+1} = (I - \underbrace{\rho_k s_k s_k^T}_{V_k^T}) H_k (I - \underbrace{\rho_k y_k y_k^T}_{V_k}) + \underbrace{\rho_k s_k s_k^T}_{0 \text{ } 1}$$

$$\begin{aligned} H_0 &= H_0 \\ H_1 &= V_0^T H_0 V_0 + \rho_0 s_0 s_0^T \\ H_2 &= V_1^T V_0^T H_0 V_0 V_1 + \rho_0 V_1^T s_0 s_0^T V_1 + \rho_1 s_1 s_1^T \\ H_3 &= V_2^T V_1^T V_0^T H_0 V_0 V_1 V_2 + \rho_0 V_2^T V_1^T s_0 s_0^T V_1 V_2 + \rho_1 V_2^T s_1 s_1^T V_2 \\ &\quad + \rho_2 s_2 s_2^T \\ &\vdots \end{aligned}$$

$$\begin{aligned} H_k &= V_{k-1}^T V_{k-2}^T \dots V_2^T H_0 V_0 V_1 \dots V_{k-1} \\ &\quad + \rho_0 V_{k-1}^T V_{k-2}^T \dots V_0^T s_0 s_0^T V_0 V_1 \dots V_{k-1} \\ &\quad + \rho_1 V_{k-1}^T \dots V_2^T s_1 s_1^T V_2 V_3 \dots V_{k-1} \\ &\quad + \rho_{k-2} V_{k-1}^T s_{k-2} s_{k-2}^T V_{k-1} \\ &\quad + \rho_{k-1} s_{k-1} s_{k-1}^T \end{aligned}$$

$q, V_{k-1} q, V_{k-2} V_{k-1} q, \dots, V_{k-1} V_{k-2} \dots V_{k-1} q$

How to compute  $H_k \cdot q$ ?

$$r = q$$

For  $i = k-1, \dots, 0$

$$\alpha_i = s_i^T r$$

$$r = V_i r$$

end

$$r = H_0 r$$

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

$$r = V_i^T r$$

$$r = r + \rho_i s_i \alpha_i$$

end

In this algo. each step:  $V_i^T r$  or  $V_i^T r$

6-8

~~$$V_i = I - \rho_i s_i s_i^T$$~~

$$V_i = I - \rho_i s_i s_i^T$$

$$V_i \cdot r = r - \rho_i s_i (s_i^T r) \Rightarrow O(n)$$

$\square \begin{pmatrix} \square \end{pmatrix}$

The total time complexity:  $O(nk) \ll O(n^2)$

$$V_i^T r = (I - \rho_i s_i s_i^T) r$$

$$= r - \rho_i s_i (s_i^T r)$$

$\square \begin{pmatrix} \square \end{pmatrix}$

L-BFGS: Only store  $m$  vectors in memory

$$\text{BFGS: } H^0 \xrightarrow{s^0, y^0} H^1 \xrightarrow{s^1, y^1} H^2 \rightarrow \dots \rightarrow H^k$$

$$\text{L-BFGS: } H_{k-m} \xrightarrow{\uparrow \substack{r \\ I}} H_{k-m+1} \xrightarrow{\uparrow \substack{r \\ 2}} \dots \xrightarrow{\uparrow \substack{r \\ 2}} H_k$$

$O(2 \cdot m \cdot n)$  memory space

$$H_k = V_{k-1}^T V_{k-2}^T \dots V_{k-m}^T H_0^k V_{k-m}^T V_{k-m+1}^T \dots V_{k-1}^T$$

$$+ \rho_{k-m}^T V_{k-1}^T \dots V_{k-m+1}^T S_{k-m} S_{k-m}^T V_{k-m+1}^T \dots V_{k-1}^T$$

$$+ \rho_{k-m+1}^T \dots$$

$$+ \rho_{k-1}^T S_{k-1} S_{k-1}^T$$



Usually. Set  $H_0^k = rI$  for some  $r > 0$ .

69

Time complexity for  $H_k \cdot p : O(mn)$

Space complexity :  $O(mn)$

Convergence  $\Rightarrow$  linear convergence.