

# Chapter 11 Quasi-Newton Methods

1. Approximating the Inverse Hessian
2. Rank One Correction Formula
3. DFP Algorithm
4. BFGS Algorithm
5. BB Algorithm



# Newton and Secant Methods

1D case:  $\min_{x \in \mathbb{R}} f(x)$

- Newton's method:

$$x^{(k+1)} = x^{(k)} - \frac{f'(x^{(k)})}{f''(x^{(k)})}.$$

- secant method:

$$x^{(k+1)} = x^{(k)} - \frac{f'(x^{(k)})}{F(x^{(k)})},$$

where

$F(x^{(k)}) = \frac{f'(x^{(k)}) - f'(x^{(k-1)})}{x^{(k)} - x^{(k-1)}}$  is  
an approximation of  $f''(x^{(k)})$ .

$n$ D case:

- Newton's method:  $\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}^{(k)},$$

where  $\mathbf{F}(\mathbf{x}^{(k)}) = \nabla^2 f(\mathbf{x}^{(k)})$ ,  
 $\mathbf{g}^{(k)} = \nabla f(\mathbf{x}^{(k)})$ .

- quasi-Newton/secant methods:  
approximates  $\mathbf{F}(\mathbf{x}^{(k)})$  or  $[\mathbf{F}(\mathbf{x}^{(k)})]^{-1}$ , i.e.,  
$$\begin{cases} \mathbf{H}_k \approx [\mathbf{F}(\mathbf{x}^{(k)})]^{-1}, \\ \mathbf{B}_k \approx \mathbf{F}(\mathbf{x}^{(k)}), \end{cases}$$

$$\text{version 1: } \mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \mathbf{H}_k \mathbf{g}^{(k)},$$

$$\text{version 2: } \mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - (\mathbf{B}_k)^{-1} \mathbf{g}^{(k)}.$$



# Version 1: $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \mathbf{H}_k \mathbf{g}^{(k)}$

## Premises on $\mathbf{H}_k$

- ①  $\mathbf{H}_k \approx [\mathbf{F}(\mathbf{x}^{(k)})]^{-1}$  and  $\mathbf{H}_k$  is symmetric positive definite;
- ② trivial computation from  $\mathbf{H}_k$  to  $\mathbf{H}_{k+1}$ , e.g.,  $\mathbf{H}_{k+1} = \mathbf{H}_k + \Delta_k$ , where  $\Delta_k$  is a “small” matrix.

**Analysis:** properties of the Hessian matrix  $\mathbf{F}(\mathbf{x}^{(k+1)})$ :

- 1D case:  $\because f'(x^{(k)}) \approx f'(x^{(k+1)}) + f''(x^{(k+1)})(x^{(k)} - x^{(k+1)})$ ,

$$\therefore f''(x^{(k+1)}) \approx \frac{f'(x^{(k+1)}) - f'(x^{(k)})}{x^{(k+1)} - x^{(k)}}$$

- nD case:  $\because \nabla f(\mathbf{x}^{(k)}) \approx \nabla f(\mathbf{x}^{(k+1)}) + \nabla^2 f(\mathbf{x}^{(k+1)})(\mathbf{x}^{(k)} - \mathbf{x}^{(k+1)})$ ,

$$\therefore \underbrace{\nabla f(\mathbf{x}^{(k)}) - \nabla f(\mathbf{x}^{(k+1)})}_{\mathbf{y}^{(k)}} \approx \underbrace{\nabla^2 f(\mathbf{x}^{(k+1)})}_{\mathbf{F}(\mathbf{x}^{(k+1)})} \underbrace{(\mathbf{x}^{(k)} - \mathbf{x}^{(k+1)})}_{\mathbf{s}^{(k)}}.$$

$$\iff \boxed{\mathbf{y}^{(k)} \approx \mathbf{F}(\mathbf{x}^{(k+1)}) \mathbf{s}^{(k)}} \iff \boxed{[\mathbf{F}(\mathbf{x}^{(k+1)})]^{-1} \mathbf{y}^{(k)} \approx \mathbf{s}^{(k)}}.$$

Replacing  $[\mathbf{F}(\mathbf{x}^{(k+1)})]^{-1}$  with  $\mathbf{H}_{k+1}$  and oblige equation to be true, we have the quasi-Newton equation, also called secant equation.



# Rank-One Correction

$$\text{quasi-Newton equation: } \mathbf{H}_{k+1} \mathbf{y}^{(k)} = \mathbf{s}^{(k)}. \quad (1)$$

- ★ solution of (1) is nonunique because there are  $\frac{n(n+1)}{2}$  degrees of freedom in symmetric matrix, and the secant equation represents  $n$  conditions. Thus, there are many choices on  $\mathbf{H}_{k+1}$ .
- ★ To guarantee  $\mathbf{H}_{k+1} \succ 0$ , it suffices the curvature condition  $\mathbf{s}^{(k)\top} \mathbf{y}^{(k)} > 0$  (which holds if the strong Wolfe conditions is used in line search).

## rank-one correction of $\mathbf{H}_{k+1}$

Given  $\mathbf{H}_k$ ,  $\mathbf{y}^{(k)}$ ,  $\mathbf{s}^{(k)}$ , update  $\mathbf{H}_{k+1}$  by solving the systems

$$\begin{cases} \mathbf{H}_{k+1} \mathbf{y}^{(k)} = \mathbf{s}^{(k)} \\ \mathbf{H}_{k+1} = \mathbf{H}_k + \mathbf{a} \mathbf{u} \mathbf{u}^\top \end{cases} \implies \dots \xRightarrow{\text{details}} \dots \implies \dots$$
$$\implies \begin{cases} \mathbf{u} = \mathbf{s}^{(k)} - \mathbf{H}_k \mathbf{y}^{(k)} \\ \mathbf{a} = \frac{1}{\mathbf{y}^{(k)\top} (\mathbf{s}^{(k)} - \mathbf{H}_k \mathbf{y}^{(k)})} \end{cases} \implies \mathbf{H}_{k+1} = \mathbf{H}_k + \frac{(\mathbf{s}^{(k)} - \mathbf{H}_k \mathbf{y}^{(k)}) (\mathbf{s}^{(k)} - \mathbf{H}_k \mathbf{y}^{(k)})^\top}{\mathbf{y}^{(k)\top} (\mathbf{s}^{(k)} - \mathbf{H}_k \mathbf{y}^{(k)})}.$$

- ★ all the premises for  $\mathbf{H}_{k+1}$  are fulfilled (check them).

# Conjugate Property

## Theorem (conjugate property)

If quasi-Newton method is applied to minimizing a quadratic function with Hessian  $\mathbf{Q} = \mathbf{Q}^\top$  such that for  $0 \leq k < n - 1$ ,  $\mathbf{H}_{k+1}\mathbf{y}^{(i)} = \mathbf{s}^{(i)}$  for all  $0 \leq i \leq k$ , If  $\alpha_i \neq 0$ ,  $0 \leq i \leq k$ , then  $\{\mathbf{d}^{(i)} = \mathbf{H}_i \mathbf{g}^{(i)}\}_{i=0}^{k+1}$  are  $\mathbf{Q}$ -conjugate.

proof. (induction): recall that  $\mathbf{s}^{(k)} = \mathbf{x}^{(k)} - \mathbf{x}^{(k+1)}$ ,

$$\therefore \mathbf{y}^{(k)} = \mathbf{g}^{(k)} - \mathbf{g}^{(k+1)} \underset{\text{function}}{\overset{\text{quadratic}}{=}} \mathbf{Q}(\mathbf{x}^{(k)} - \mathbf{x}^{(k+1)}) = \mathbf{Q}\mathbf{s}^{(k)},$$

Recall that  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \mathbf{H}_k \mathbf{g}^{(k)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)}$ .

(1) case of  $k = 0$ .  $\therefore \alpha_0 \neq 0$ ,  $\therefore \mathbf{d}^{(0)} = -\frac{\mathbf{s}^{(0)}}{\alpha_0}$ .

$$\begin{aligned} \therefore (\mathbf{d}^{(1)})^\top \mathbf{Q} \mathbf{d}^{(0)} &= -(\mathbf{g}^{(1)})^\top \mathbf{H}_1 \mathbf{Q} \mathbf{d}^{(0)} = (\mathbf{g}^{(1)})^\top \mathbf{H}_1 \frac{\mathbf{Q} \mathbf{s}^{(0)}}{\alpha_0} \\ &= (\mathbf{g}^{(1)})^\top \frac{\mathbf{H}_1 \mathbf{y}^{(0)}}{\alpha_0} \underset{\text{equation}}{\overset{\text{secant}}{=}} (\mathbf{g}^{(1)})^\top \frac{\mathbf{s}^{(0)}}{\alpha_0} = (\mathbf{g}^{(1)})^\top \mathbf{d}^{(0)} \underset{\text{exact}}{\overset{\alpha_k}{=}} 0, \end{aligned}$$

$\therefore \mathbf{d}^{(0)}$  and  $\mathbf{d}^{(1)}$  are  $\mathbf{Q}$ -conjugate.



# Conjugate Property

- (2) Assume the conjugacy holds for  $k - 1$  (where  $k < n - 1$ ). We now prove the the conjugate statement hold for  $k$ , i.e.,  $\{\mathbf{d}^{(k+1)}\}_{i=0}^{k+1}$  are  $\mathbf{Q}$ -conjugate. It suffices to show that

$$(\mathbf{d}^{(k+1)})^\top \mathbf{Q} \mathbf{d}^{(i)} = 0, \quad 0 \leq i \leq k.$$

Given  $i$  ( $0 \leq i \leq k$ ), using the same algebraic steps as  $k = 0$ , and assumption that  $\alpha_i \neq 0$ , we have

$$(\mathbf{d}^{(k+1)})^\top \mathbf{Q} \mathbf{d}^{(i)} = -(\mathbf{g}^{(k+1)})^\top \mathbf{H}_{k+1} \mathbf{Q} \mathbf{d}^{(i)} = \dots = -(\mathbf{g}^{(k+1)})^\top \mathbf{d}^{(i)}$$

$\therefore \mathbf{d}^{(0)}, \dots, \mathbf{d}^{(k)}$  are  $\mathbf{Q}$ -conjugate by assumption, we have

$$(\mathbf{g}^{(k+1)})^\top \mathbf{d}^{(i)} = 0.$$

$$\therefore (\mathbf{d}^{(k+1)})^\top \mathbf{Q} \mathbf{d}^{(i)} = 0.$$

★ quasi-Newton method solves a quadratic of  $n$  variables in at most  $n$  steps.



# Rank-One Quasi-Newton Method

**Require:**  $\mathbf{x}^{(0)} \in \mathbb{R}^n$  and  $\mathbf{H}_0 \succ 0$ . set  $k := 0$ ;

1: **repeat**

2:  $\mathbf{d}^{(k)} = -\mathbf{H}_k \mathbf{g}^{(k)} \implies$  quasi-Newton direction

3:  $\alpha_k = \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)}) \implies$  compute step size

4:  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{d}^{(k)} \implies$  updating iterative point.

5:  $\mathbf{s}^{(k)} = \mathbf{x}^{(k)} - \mathbf{x}^{(k+1)}, \quad \mathbf{y}^{(k)} = \mathbf{g}^{(k)} - \mathbf{g}^{(k+1)},$

6:  $\mathbf{H}_{k+1} = \mathbf{H}_k + \frac{(\mathbf{s}^{(k)} - \mathbf{H}_k \mathbf{y}^{(k)})(\mathbf{s}^{(k)} - \mathbf{H}_k \mathbf{y}^{(k)})^\top}{\mathbf{y}^{(k)\top}(\mathbf{s}^{(k)} - \mathbf{H}_k \mathbf{y}^{(k)})} \implies$  updating  $\mathbf{H}_k$ ;

7: **until**  $\|\mathbf{g}^{(k)}\| \leq \epsilon$ .

- ★  $\mathbf{H}_{k+1}$  generated by rank-one algorithm may be  $\mathbf{H}_{k+1} \not\succ 0$  and thus  $\mathbf{d}^{(k+1)}$  may not be a descent direction. This happens even in the quadratic case.
- ★ if  $(\mathbf{y}^{(k)})^\top (\mathbf{s}^{(k)} - \mathbf{H}_k \mathbf{y}^{(k)})$  is close to zero, then there may be numerical problems in evaluating  $\mathbf{H}_{k+1}$ .

An example of  $\mathbf{H}_{k+1} \not\succ 0$ : see 11.2 in textbook.

# Rank-One Algorithm

Example (apply rank-one algorithm to minimize  $f$ )

$$f(x_1, x_2) = x_1^2 + \frac{1}{2}x_2^2 + 3, \mathbf{x}^{(0)} = [1, 2]^\top \text{ and } \mathbf{H}_0 = \mathbf{I}_2.$$

Ans:  $\mathbf{g}^{(k)} = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \mathbf{x}^{(k)}.$

S1:  $\mathbf{d}^{(0)} = -\mathbf{g}^{(0)} = [-2, -2]^\top,$

$$\alpha_0 = \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(0)} + \alpha \mathbf{d}^{(0)}) = -\frac{\mathbf{g}^{(0)\top} \mathbf{d}^{(0)}}{\mathbf{d}^{(0)\top} \mathbf{Q} \mathbf{d}^{(0)}} = \frac{2}{3},$$

$$\mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \alpha_0 \mathbf{d}^{(0)} = \left[-\frac{1}{3}, \frac{2}{3}\right]^\top.$$

S2:  $\mathbf{g}^{(1)} = \mathbf{Q} \mathbf{x}^{(1)} = \left[-\frac{2}{3}, \frac{2}{3}\right]^\top,$

$$\mathbf{s}^{(0)} = \mathbf{x}^{(1)} - \mathbf{x}^{(0)} = \left[-\frac{4}{3}, -\frac{4}{3}\right]^\top, \mathbf{y}^{(0)} = \mathbf{g}^{(1)} - \mathbf{g}^{(0)} = \left[-\frac{8}{3}, -\frac{4}{3}\right]^\top,$$

$$\mathbf{H}_1 = \mathbf{H}_0 + \frac{(\mathbf{s}^{(0)} - \mathbf{H}_0 \mathbf{y}^{(0)})(\mathbf{s}^{(0)} - \mathbf{H}_0 \mathbf{y}^{(0)})^\top}{\mathbf{y}^{(0)\top} (\mathbf{s}^{(0)} - \mathbf{H}_0 \mathbf{y}^{(0)})} = \begin{bmatrix} \frac{1}{2} & 0 \\ 0 & 1 \end{bmatrix},$$

$$\mathbf{d}^{(1)} = -\mathbf{H}_1 \mathbf{g}^{(1)} = \left[\frac{1}{3}, -\frac{2}{3}\right]^\top,$$

$$\alpha_1 = -\frac{\mathbf{g}^{(1)\top} \mathbf{d}^{(1)}}{\mathbf{d}^{(1)\top} \mathbf{Q} \mathbf{d}^{(1)}} = 1, \mathbf{x}^{(2)} = \mathbf{x}^{(1)} + \alpha_1 \mathbf{d}^{(1)} = [0, 0]^\top,$$

$$\therefore \mathbf{g}^{(2)} = \mathbf{0}, \therefore \mathbf{x}^* = \mathbf{x}^{(2)}.$$





# DFP Algorithm

## quasi-Newton equation

$$\mathbf{H}_{k+1} \mathbf{y}^{(k)} = \mathbf{s}^{(k)}, \text{ where } \mathbf{s}^{(k)} = \mathbf{x}^{(k)} - \mathbf{x}^{(k+1)}, \mathbf{y}^{(k)} = \mathbf{g}^{(k)} - \mathbf{g}^{(k+1)}$$

## rank-two correction of $\mathbf{H}_{k+1}$ : DFP formula

Given  $\mathbf{H}_k$ ,  $\mathbf{y}^{(k)}$ ,  $\mathbf{s}^{(k)}$ , update  $\mathbf{H}_{k+1}$  by solving the systems

$$\begin{cases} \mathbf{H}_{k+1} \mathbf{y}^{(k)} = \mathbf{s}^{(k)} \\ \mathbf{H}_{k+1} = \mathbf{H}_k + \mathbf{a} \mathbf{u} \mathbf{u}^\top + \mathbf{b} \mathbf{v} \mathbf{v}^\top \end{cases} \xRightarrow{\text{details}} \dots \Rightarrow$$



# DFP Algorithm

## quasi-Newton equation

$$\mathbf{H}_{k+1} \mathbf{y}^{(k)} = \mathbf{s}^{(k)}, \text{ where } \mathbf{s}^{(k)} = \mathbf{x}^{(k)} - \mathbf{x}^{(k+1)}, \mathbf{y}^{(k)} = \mathbf{g}^{(k)} - \mathbf{g}^{(k+1)}$$

## rank-two correction of $\mathbf{H}_{k+1}$ : DFP formula

Given  $\mathbf{H}_k$ ,  $\mathbf{y}^{(k)}$ ,  $\mathbf{s}^{(k)}$ , update  $\mathbf{H}_{k+1}$  by solving the systems

$$\begin{cases} \mathbf{H}_{k+1} \mathbf{y}^{(k)} = \mathbf{s}^{(k)} \\ \mathbf{H}_{k+1} = \mathbf{H}_k + a \mathbf{u} \mathbf{u}^\top + b \mathbf{v} \mathbf{v}^\top \end{cases} \xRightarrow{\text{details}} \dots \Rightarrow \begin{cases} \mathbf{u} = \mathbf{s}^{(k)}, \\ \mathbf{v} = \mathbf{H}_k \mathbf{y}^{(k)}, \\ a = \frac{1}{\mathbf{s}^{(k)\top} \mathbf{y}^{(k)}}, \\ b = -\frac{1}{\mathbf{y}^{(k)\top} \mathbf{H}_k \mathbf{y}^{(k)}}. \end{cases}$$

$$\Rightarrow \mathbf{H}_{k+1} = \mathbf{H}_k + \frac{\mathbf{s}^{(k)} \mathbf{s}^{(k)\top}}{\mathbf{s}^{(k)\top} \mathbf{y}^{(k)}} - \frac{[\mathbf{H}_k \mathbf{y}^{(k)}][\mathbf{H}_k \mathbf{y}^{(k)}]^\top}{\mathbf{y}^{(k)\top} \mathbf{H}_k \mathbf{y}^{(k)}}.$$



# DFP Algorithm

Example (apply DFP algorithm to minimize  $f$ )

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top \begin{bmatrix} 4 & 2 \\ 2 & 2 \end{bmatrix} \mathbf{x} - \mathbf{x}^\top \begin{bmatrix} -1 \\ 1 \end{bmatrix}. \text{ The initial point } \mathbf{x}^{(0)} = [0, 0]^\top \text{ and } \mathbf{H}_0 = \mathbf{I}_2.$$

Ans:  $\mathbf{g}^{(k)} = \begin{bmatrix} 4 & 2 \\ 2 & 2 \end{bmatrix} \mathbf{x}^{(k)} - \begin{bmatrix} -1 \\ 1 \end{bmatrix}.$

S1:  $\mathbf{g}^{(0)} = [1, -1]^\top$ ,  $\mathbf{d}^{(0)} = -\mathbf{H}_0 \mathbf{g}^{(0)} = [-1, 1]^\top$ ,  
 $\alpha_0 = \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(0)} + \alpha \mathbf{d}^{(0)}) = 1$ ,  $\mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \alpha_0 \mathbf{d}^{(0)} = [-1, 1]^\top$ .

S2:  $\mathbf{g}^{(1)} = [-1, -1]^\top$ ,  
 $\mathbf{y}^{(0)} = \mathbf{g}^{(1)} - \mathbf{g}^{(0)} = [-2, 0]^\top$ ,  $\mathbf{s}^{(0)} = \mathbf{x}^{(1)} - \mathbf{x}^{(0)} = [-1, 1]^\top$ ,

$$\mathbf{H}_1 = \mathbf{H}_0 + \frac{\mathbf{s}^{(0)} \mathbf{s}^{(0)\top}}{\mathbf{s}^{(0)\top} \mathbf{y}^{(0)}} - \frac{(\mathbf{H}_0 \mathbf{y}^{(0)}) (\mathbf{H}_0 \mathbf{y}^{(0)})^\top}{\mathbf{y}^{(0)\top} \mathbf{H}_0 \mathbf{y}^{(0)}} = \begin{bmatrix} \frac{1}{2} & -\frac{1}{2} \\ -\frac{1}{2} & \frac{3}{2} \end{bmatrix},$$

$$\mathbf{d}^{(1)} = -\mathbf{H}_1 \mathbf{g}^{(1)} = [0, 1]^\top,$$

$$\alpha_1 = \arg \min_{\alpha \geq 0} f(\mathbf{x}^{(1)} + \alpha \mathbf{d}^{(1)}) = \frac{1}{2}, \quad \mathbf{x}^{(2)} = \mathbf{x}^{(1)} + \alpha_1 \mathbf{d}^{(1)} = [-1, \frac{3}{2}]^\top.$$

$\therefore \mathbf{g}^{(2)} = [0, 0]^\top$ ,  $\therefore \mathbf{x}^{(2)}$  is the minimizer.

★  $\mathbf{d}^{(0)}$  and  $\mathbf{d}^{(1)}$  are  $Q$ -conjugate (check it).



# DFP Algorithm

Example (apply DFP algorithm to minimize  $f$ )

$$f(\mathbf{x}) = 2x_1^2 + x_2^2 - 4x_1 + 2 \text{ with } \mathbf{H}_0 = \mathbf{I}_2 \text{ and } \mathbf{x}^{(0)} = [2, 1]^\top.$$

Ans:  $\nabla f(\mathbf{x}) = \begin{bmatrix} 4x_1 - 4 \\ 2x_2 \end{bmatrix}$ ,  $\nabla^2 f(\mathbf{x}) = \begin{bmatrix} 4 & 0 \\ 0 & 2 \end{bmatrix} =: \mathbf{Q}$ ,

S1:  $\mathbf{g}^{(0)} = [4, 2]^\top$ ,  $\mathbf{d}_0 = -\mathbf{H}_0 \mathbf{g}^{(0)} = [-4, -2]^\top$ ,  
 $\alpha_0 = -\frac{\mathbf{g}^{(0)\top} \mathbf{d}^{(0)}}{\mathbf{d}^{(0)\top} \mathbf{Q} \mathbf{d}^{(0)}} = \frac{5}{18}$ ,  $\mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \alpha_0 \mathbf{d}^{(0)} = [\frac{8}{9}, \frac{4}{9}]^\top$ .

S2:  $\mathbf{s}^{(0)} = \mathbf{x}^{(1)} - \mathbf{x}^{(0)} = [-\frac{10}{9}, -\frac{5}{9}]^\top$ ,  $\mathbf{g}^{(1)} = [-\frac{4}{9}, \frac{8}{9}]^\top$ ,  
 $\mathbf{y}^{(0)} = \mathbf{g}^{(1)} - \mathbf{g}^{(0)} = [-\frac{40}{9}, -\frac{10}{9}]^\top$ ,

$$\mathbf{H}_1 = \mathbf{H}_0 + \frac{\mathbf{s}^{(0)} \mathbf{s}^{(0)\top}}{\mathbf{s}^{(0)\top} \mathbf{y}^{(0)}} - \frac{\mathbf{H}_0 \mathbf{y}^{(0)} \mathbf{y}^{(0)\top} \mathbf{H}_0}{\mathbf{y}^{(0)\top} \mathbf{H}_0 \mathbf{y}^{(0)}} = \frac{1}{306} \begin{bmatrix} 86 & -38 \\ -38 & 305 \end{bmatrix},$$

$$\mathbf{d}^{(1)} = -\mathbf{H}_1 \mathbf{g}^{(1)} = [\frac{4}{17}, -\frac{16}{17}]^\top,$$

$$\alpha_1 = -\frac{\mathbf{g}^{(1)\top} \mathbf{d}^{(1)}}{\mathbf{d}^{(1)\top} \mathbf{Q} \mathbf{d}^{(1)}} = \frac{17}{36},$$

$$\mathbf{x}^{(2)} = \mathbf{x}^{(1)} + \alpha_1 \mathbf{d}^{(1)} = [1, 0]^\top,$$

$\therefore \mathbf{g}^{(2)} = \mathbf{0}$ ,  $\therefore \mathbf{x}^{(2)}$  is the minimizer.

★  $\mathbf{d}^{(0)}$  and  $\mathbf{d}^{(1)}$  are  $\mathbf{Q}$ -conjugate (check it).



## Theorem (inheritance of DFP)

*Suppose that  $\mathbf{g}^{(k)} \neq \mathbf{0}$ . In the DFP algorithm, if  $\mathbf{H}_k \succ 0$ , then so is  $\mathbf{H}_{k+1}$ .*

- ★ DFP was developed by Davidon in 1959 and was modified by Fletcher and Powell in 1963.
- ★ DFP is superior to rank-one algorithm as it can guarantee  $\mathbf{H}_k \succ 0$ .
- ★ DFP may get stuck for large nonlinear optimization because  $\mathbf{H}_k$  is close to singular.



# BFGS Algorithm

**History:** In 1970, an alternative update formula was suggested independently by Broyden, Fletcher, Goldfarb, and Shanno, which is called BFGS algorithm.

## motivation

Newton's method:  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \mathbf{F}(\mathbf{x}^{(k)})^{-1} \mathbf{g}^{(k)}$ .

Instead of seeking  $\mathbf{H}_k \approx [\mathbf{F}(\mathbf{x}^{(k)})]^{-1}$ , we alternatively find an approximation of Hessian, i.e.,  $\mathbf{B}_k \approx \mathbf{F}(\mathbf{x}^{(k)})$ , thus

quasi-Newton methods:  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \mathbf{B}_k^{-1} \mathbf{g}^{(k)}$ .

## choice of $\mathbf{B}_k$ (analogously to $\mathbf{H}_k$ )

$$\begin{aligned}\therefore \nabla f(\mathbf{x}^{(k)}) &\approx \nabla f(\mathbf{x}^{(k+1)}) + \mathbf{F}(\mathbf{x}^{(k+1)})(\mathbf{x}^{(k)} - \mathbf{x}^{(k+1)}), \\ \therefore \underbrace{\nabla f(\mathbf{x}^{(k)}) - \nabla f(\mathbf{x}^{(k+1)})}_{\mathbf{y}^{(k)}} &\approx \mathbf{F}(\mathbf{x}^{(k+1)}) \underbrace{(\mathbf{x}^{(k)} - \mathbf{x}^{(k+1)})}_{\mathbf{s}^{(k)}}.\end{aligned}$$

$$\therefore \mathbf{y}^{(k)} \approx \mathbf{F}(\mathbf{x}^{(k+1)}) \mathbf{s}^{(k)}.$$

$$\therefore \text{Quasi-Newton equation: } \mathbf{y}^{(k)} = \mathbf{B}_{k+1} \mathbf{s}^{(k)}.$$

# BFGS Algorithm

## quasi-Newton equation

$$\mathbf{B}_{k+1} \mathbf{s}^{(k)} = \mathbf{y}^{(k)}, \text{ where } \mathbf{s}^{(k)} = \mathbf{x}^{(k)} - \mathbf{x}^{(k+1)}, \mathbf{y}^{(k)} = \mathbf{g}^{(k)} - \mathbf{g}^{(k+1)}$$

## rank-two correction of $\mathbf{B}_{k+1}$ : BFGS formula

Given  $\mathbf{B}_k$ ,  $\mathbf{y}^{(k)}$ ,  $\mathbf{s}^{(k)}$ , update  $\mathbf{B}_{k+1}$  by solving the systems

$$\begin{cases} \mathbf{B}_{k+1} \mathbf{s}^{(k)} = \mathbf{y}^{(k)}, \\ \mathbf{B}_{k+1} = \mathbf{B}_k + \mathbf{a} \mathbf{u} \mathbf{u}^\top + \mathbf{b} \mathbf{v} \mathbf{v}^\top, \end{cases} \quad \Rightarrow \dots \xRightarrow{\text{detail}} \dots$$



# BFGS Algorithm

## quasi-Newton equation

$$\mathbf{B}_{k+1} \mathbf{s}^{(k)} = \mathbf{y}^{(k)}, \text{ where } \mathbf{s}^{(k)} = \mathbf{x}^{(k)} - \mathbf{x}^{(k+1)}, \mathbf{y}^{(k)} = \mathbf{g}^{(k)} - \mathbf{g}^{(k+1)}$$

## rank-two correction of $\mathbf{B}_{k+1}$ : BFGS formula

Given  $\mathbf{B}_k$ ,  $\mathbf{y}^{(k)}$ ,  $\mathbf{s}^{(k)}$ , update  $\mathbf{B}_{k+1}$  by solving the systems

$$\begin{cases} \mathbf{B}_{k+1} \mathbf{s}^{(k)} = \mathbf{y}^{(k)}, \\ \mathbf{B}_{k+1} = \mathbf{B}_k + \mathbf{a} \mathbf{u} \mathbf{u}^\top + \mathbf{b} \mathbf{v} \mathbf{v}^\top, \end{cases} \implies \dots \xRightarrow{\text{detail}} \dots \begin{cases} \mathbf{u} = \mathbf{y}^{(k)}, \\ \mathbf{v} = \mathbf{B}_k \mathbf{s}^{(k)}, \\ \mathbf{a} = \frac{1}{\mathbf{s}^{(k)\top} \mathbf{y}^{(k)}}, \\ \mathbf{b} = -\frac{1}{\mathbf{s}^{(k)\top} \mathbf{B}_k \mathbf{s}^{(k)}}. \end{cases}$$
$$\implies \mathbf{B}_{k+1} = \mathbf{B}_k + \frac{\mathbf{y}^{(k)} \mathbf{y}^{(k)\top}}{\mathbf{y}^{(k)\top} \mathbf{s}^{(k)}} - \frac{\mathbf{B}_k \mathbf{s}^{(k)} \mathbf{s}^{(k)\top} \mathbf{B}_k}{\mathbf{s}^{(k)\top} \mathbf{B}_k \mathbf{s}^{(k)}}.$$

★ Recall DFP formula  $\mathbf{H}_{k+1}^{\text{DFP}} = \mathbf{H}_k + \frac{\mathbf{s}^{(k)} \mathbf{s}^{(k)\top}}{\mathbf{s}^{(k)\top} \mathbf{y}^{(k)}} - \frac{\mathbf{H}_k \mathbf{y}^{(k)} \mathbf{y}^{(k)\top} \mathbf{H}_k}{\mathbf{y}^{(k)\top} \mathbf{H}_k \mathbf{y}^{(k)}}$ , we have  $\mathbf{H}_{k+1}$  and  $\mathbf{B}_{k+1}$  are “dual” of each other, i.e.,  $\mathbf{B}_{k+1}$  can be derived by commuting  $\mathbf{H}_k \leftrightarrow \mathbf{B}_k$ ,  $\mathbf{y}^{(k)} \leftrightarrow \mathbf{s}^{(k)}$ .





# BFGS Algorithm

quasi-Newton methods with BFGS:  $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \mathbf{B}_k^{-1} \mathbf{g}^{(k)}$ .

**Question:** how to compute the inverse? any skill to derive  $\mathbf{B}_{k+1}^{-1}$  from  $\mathbf{B}_k^{-1}$ ?

$$\mathbf{B}_{k+1}^{-1} = \left( \mathbf{B}_k + \frac{\mathbf{y}^{(k)} \mathbf{y}^{(k)\top}}{\mathbf{y}^{(k)\top} \mathbf{s}^{(k)}} - \frac{\mathbf{B}_k \mathbf{s}^{(k)} \mathbf{s}^{(k)\top} \mathbf{B}_k}{\mathbf{s}^{(k)\top} \mathbf{B}_k \mathbf{s}^{(k)}} \right)^{-1} = \boxed{?}$$

## Lemma (Sherman-Morrison formula)

Let  $\mathbf{U} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{V} \in \mathbb{R}^{n \times m}$ . If  $\mathbf{I}_m + \mathbf{U}\mathbf{V}$  is invertible, then  $\mathbf{I}_n + \mathbf{V}\mathbf{U}$  is also invertible and  $(\mathbf{I}_n + \mathbf{V}\mathbf{U})^{-1} = \mathbf{I}_n - \mathbf{V}(\mathbf{I}_m + \mathbf{U}\mathbf{V})^{-1}\mathbf{U}$ .

**proof.**  $\mathbf{V} + \mathbf{V}\mathbf{U}\mathbf{V} = \mathbf{V}(\mathbf{I}_m + \mathbf{U}\mathbf{V}) = (\mathbf{I}_n + \mathbf{V}\mathbf{U})\mathbf{V}$ .

$$\therefore \mathbf{V} = (\mathbf{I}_n + \mathbf{V}\mathbf{U})\mathbf{V}(\mathbf{I}_m + \mathbf{U}\mathbf{V})^{-1}.$$

$$\begin{aligned} \therefore \mathbf{I}_n &= \mathbf{I}_n + \mathbf{V}\mathbf{U} - \mathbf{V}\mathbf{U} = \mathbf{I}_n + \mathbf{V}\mathbf{U} - (\mathbf{I}_n + \mathbf{V}\mathbf{U})\mathbf{V}(\mathbf{I}_m + \mathbf{U}\mathbf{V})^{-1}\mathbf{U} \\ &= (\mathbf{I}_n + \mathbf{V}\mathbf{U})[\mathbf{I}_n - \mathbf{V}(\mathbf{I}_m + \mathbf{U}\mathbf{V})^{-1}\mathbf{U}] \end{aligned}$$

**Corollary** ( $\mathbf{A} \in \mathbb{R}^{m \times m}$  be invertible,  $\mathbf{U} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{V} \in \mathbb{R}^{n \times m}$ ,  $\mathbf{u}$  and  $\mathbf{v} \in \mathbb{R}^m$ )

$$(\mathbf{A} + \mathbf{U}\mathbf{V})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{U}(\mathbf{I}_n + \mathbf{V}\mathbf{A}^{-1}\mathbf{U})^{-1}\mathbf{V}\mathbf{A}^{-1}.$$

$$(\mathbf{A} + \mathbf{u}\mathbf{v}^\top)^{-1} = \mathbf{A}^{-1} - \frac{(\mathbf{A}^{-1}\mathbf{u})(\mathbf{v}^\top \mathbf{A}^{-1})}{1 + \mathbf{v}^\top \mathbf{A}^{-1}\mathbf{u}}. \text{ (proof. handwriting the details)}$$

# BFGS Algorithm

inverse of  $B_{k+1}$

Assume  $B_k^{-1}$  is given (denoted by  $H_k$ ). By applying S-M formula twice, we have

$$\begin{aligned} B_{k+1}^{-1} &= H_{k+1} \\ &= H_k + \left(1 + \frac{\mathbf{y}^{(k)\top} H_k \mathbf{y}^{(k)}}{\mathbf{y}^{(k)\top} \mathbf{s}^{(k)}}\right) \frac{\mathbf{s}^{(k)} \mathbf{s}^{(k)\top}}{\mathbf{s}^{(k)\top} \mathbf{y}^{(k)}} - \frac{H_k \mathbf{y}^{(k)} \mathbf{s}^{(k)\top} + (H_k \mathbf{y}^{(k)} \mathbf{s}^{(k)\top})^\top}{\mathbf{y}^{(k)\top} \mathbf{s}^{(k)}}, \end{aligned}$$

## discussion on DFP and BFGS

- Both are rank-two formulae.

- DFP:  $\mathbf{H}_{k+1} \approx [\mathbf{F}(\mathbf{x}^{(k+1)})]^{-1}$ ,  $\mathbf{H}_{k+1} = \mathbf{H}_k + \frac{\mathbf{s}^{(k)} \mathbf{s}^{(k)\top}}{\mathbf{s}^{(k)\top} \mathbf{y}^{(k)}} - \frac{\mathbf{H}_k \mathbf{y}^{(k)} \mathbf{y}^{(k)\top} \mathbf{H}_k}{\mathbf{y}^{(k)\top} \mathbf{H}_k \mathbf{y}^{(k)}}.$

- BFGS:  $\mathbf{B}_{k+1} \approx \mathbf{F}(\mathbf{x}^{(k+1)})$ ,  $\mathbf{B}_{k+1} = \mathbf{B}_k + \frac{\mathbf{y}^{(k)} \mathbf{y}^{(k)\top}}{\mathbf{y}^{(k)\top} \mathbf{s}^{(k)}} - \frac{\mathbf{B}_k \mathbf{s}^{(k)} \mathbf{s}^{(k)\top} \mathbf{B}_k}{\mathbf{s}^{(k)\top} \mathbf{B}_k \mathbf{s}^{(k)}}.$

Question:  $\mathbf{H}_{k+1} = \mathbf{B}_{k+1}^{-1}$ ?

- Combinations of DFP and BFGS:

$$\bar{\mathbf{H}}_{k+1} = \alpha \mathbf{H}_{k+1} + (1 - \alpha) \mathbf{B}_{k+1}^{-1} \quad \text{or} \quad \bar{\mathbf{B}}_{k+1} = \alpha \mathbf{B}_{k+1} + (1 - \alpha) \mathbf{H}_{k+1}^{-1}$$

- ★ BFGS fulfills all premises of quasi-Newton method, e.g., positive definiteness, symmetry when minimizing quadratic problem.



## Example (apply BFGS algorithm to minimize $f$ )

$$f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top \begin{bmatrix} 5 & -3 \\ -3 & 2 \end{bmatrix} \mathbf{x} - \mathbf{x}^\top \begin{bmatrix} 0 \\ 1 \end{bmatrix} + \log(\pi) \text{ with } \mathbf{H}_0 = \mathbf{I}_2 \text{ and } \mathbf{x}^{(0)} = [0, 0]^\top.$$

$$\text{S1: } \mathbf{d}^{(0)} = -\mathbf{g}^{(0)} = -(\mathbf{Q}\mathbf{x}^{(0)} - \mathbf{b}),$$

$$\alpha_0 = -\frac{\mathbf{g}^{(0)\top} \mathbf{d}^{(0)}}{\mathbf{d}^{(0)\top} \mathbf{Q} \mathbf{d}^{(0)}} = \frac{1}{2}, \quad \mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \alpha_0 \mathbf{d}^{(0)} = [0, \frac{1}{2}]^\top.$$

$$\text{S2: } \mathbf{s}^{(0)} = \mathbf{x}^{(1)} - \mathbf{x}^{(0)} = [0, \frac{1}{2}]^\top, \quad \mathbf{g}^{(1)} = \mathbf{Q}\mathbf{x}^{(1)} - \mathbf{b} = [\frac{-3}{2}, 0]^\top,$$

$$\mathbf{y}^{(0)} = \mathbf{g}^{(1)} - \mathbf{g}^{(0)} = [\frac{-3}{2}, 1]^\top,$$

$$\mathbf{H}_1 = \mathbf{H}_0 + \left(1 + \frac{\mathbf{y}^{(0)\top} \mathbf{H}_0 \mathbf{y}^{(0)}}{\mathbf{y}^{(0)\top} \mathbf{s}^{(0)}}\right) \frac{\mathbf{s}^{(0)} \mathbf{s}^{(0)\top}}{\mathbf{s}^{(0)\top} \mathbf{y}^{(0)}} - \frac{\mathbf{s}^{(0)} \mathbf{y}^{(0)\top} \mathbf{H}_0 + \mathbf{H}_0 \mathbf{y}^{(0)} \mathbf{s}^{(0)\top}}{\mathbf{y}^{(0)\top} \mathbf{s}^{(0)}} = \begin{bmatrix} 1 & \frac{3}{2} \\ \frac{3}{2} & \frac{11}{4} \end{bmatrix},$$

$$\mathbf{d}^{(1)} = -\mathbf{H}_1 \mathbf{g}^{(1)} = [\frac{3}{2}, \frac{9}{4}]^\top,$$

$$\alpha_1 = -\frac{\mathbf{g}^{(1)\top} \mathbf{d}^{(1)}}{\mathbf{d}^{(1)\top} \mathbf{Q} \mathbf{d}^{(1)}} = 2,$$

$$\mathbf{x}^{(2)} = \mathbf{x}^{(1)} + \alpha_1 \mathbf{d}^{(1)} = [3, 5]^\top,$$

$$\therefore \mathbf{g}^{(2)} = \mathbf{0}, \therefore \mathbf{x}^{(2)} \text{ is the minimizer.}$$



# Comparisons of Quasi-Newton and Newton Methods

Quasi-Newton method	Newton method
only need function value and gradient	need function value, gradient and Hessian
$\mathbf{H}_k \succ 0 / \mathbf{B}_k \succ 0$ for several updates	$\mathbf{F}(\mathbf{x}^{(k)})$ may be not positive definite
$O(n^2)$ multiplications in each iteration	$O(n^3)$ multiplications in each iteration

## Example (apply BFGS algorithm to minimize nonquadratic $f$ )

$f(\mathbf{x}) = f(x_1, x_2) = 100(x_2 - x_1^2)^2 + (1 - x_1)^2$  with  $\mathbf{x}^{(0)} = [-1.2, 1]^\top$  (note: exact minimizer  $\mathbf{x}^* = [1, 1]^\top$ )

steepest descent	BFGS	Newton
1.827e-4	1.70e-3	3.48e-2
1.826e-4	1.17e-3	1.44e-2
1.824e-4	1.34e-4	1.82e-4
1.823e-4	1.01e-6	1.17e-8

After coding on Matlab, Table lists some values of  $\|\mathbf{x}^{(k)} - \mathbf{x}^*\|_2$ . To reach  $\|\mathbf{g}^{(k)}\|_2 \leq 10^{-5}$ , steepest descent needs 5264 iterations, BFGS needs 34 iterations, Newton needs 21 iterations.

★ BFGS has the superlinear rate of convergence on practical problems.



# BFGS Algorithm

- ★ For nonquadratic problems, quasi-Newton algorithms may not converge in  $n$  steps.
- ★ **Restart technique:** As in the case of the conjugate gradient methods, we may reinitialize the direction vector to the negative gradient after every few iterations (e.g.,  $n$  or  $n + 1$ ), and continue until the algorithm satisfies the stopping criterion.



# Barzilai-Borwein (BB) Algorithm

## quasi-Newton method with BFGS

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \mathbf{B}_k^{-1} \mathbf{g}^{(k)}, \text{ where } \mathbf{B}_k \approx \mathbf{F}(\mathbf{x}^{(k)}).$$

**Question:** instead of rank-one/two  $\approx \mathbf{F}(\mathbf{x}^{(k)})$ , how about  $\alpha_k \mathbf{I} \approx \mathbf{F}(\mathbf{x}^{(k)})$ ?

## BB algorithm

Given  $\mathbf{y}^{(k)}, \mathbf{s}^{(k)}$ , solve systems

$$\text{BFGS} \begin{cases} \mathbf{B}_{k+1} \mathbf{s}^{(k)} = \mathbf{y}^{(k)} \\ \mathbf{B}_{k+1} = \mathbf{B}_k + a \mathbf{u} \mathbf{u}^\top + b \mathbf{v} \mathbf{v}^\top \end{cases} \xrightarrow{\text{new idea}} \begin{cases} \mathbf{B}_{k+1} \mathbf{s}^{(k)} = \mathbf{y}^{(k)} \\ \mathbf{B}_{k+1} = t_{k+1} \mathbf{I} \end{cases} \quad (*)$$

☹ (\*) has no exact solution.

😊 (\*) has a least squares solution as  $t_{k+1} = \frac{(\mathbf{y}^{(k)})^\top \mathbf{s}^{(k)}}{(\mathbf{s}^{(k)})^\top \mathbf{s}^{(k)}} \text{ (why?)}$

## iterative scheme of BB

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k^{\text{BB}} \mathbf{g}^{(k)}, \text{ where } \alpha_k^{\text{BB}} = \frac{1}{t_k} = \frac{(\mathbf{s}^{(k-1)})^\top \mathbf{s}^{(k-1)}}{(\mathbf{y}^{(k-1)})^\top \mathbf{s}^{(k-1)}}.$$

# Barzilai-Borwein (BB) Algorithm

**Question:** what is the relationship of BB and steepest descent (SD) methods when minimizing quadratic function  $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top \mathbf{Q}\mathbf{x} + \mathbf{b}^\top \mathbf{x}$ ?

## relationship of BB and SD

exact stepsize of SD method:  $\alpha_k^{\text{SD}} = \frac{(\mathbf{g}^{(k)})^\top \mathbf{g}^{(k)}}{(\mathbf{g}^{(k)})^\top \mathbf{Q} \mathbf{g}^{(k)}}.$

$$\therefore \mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \mathbf{g}^{(k)}$$

$$\therefore \mathbf{g}^{(k)} = \frac{\mathbf{x}^{(k)} - \mathbf{x}^{(k+1)}}{\alpha_k} = \frac{\mathbf{s}^{(k)}}{\alpha_k},$$

$$\therefore \alpha_k^{\text{SD}} = \frac{(\mathbf{g}^{(k)})^\top \mathbf{g}^{(k)}}{(\mathbf{g}^{(k)})^\top \mathbf{Q} \mathbf{g}^{(k)}} = \frac{(\mathbf{s}^{(k)})^\top \mathbf{s}^{(k)}}{(\mathbf{s}^{(k)})^\top \mathbf{Q} \mathbf{s}^{(k)}},$$

$$\therefore \mathbf{Q} \mathbf{s}^{(k)} = \mathbf{Q}(\mathbf{x}^{(k)} - \mathbf{x}^{(k+1)}) = \mathbf{g}^{(k)} - \mathbf{g}^{(k+1)} = \mathbf{y}^{(k)}$$

$$\therefore \alpha_k^{\text{SD}} = \frac{(\mathbf{s}^{(k)})^\top \mathbf{s}^{(k)}}{(\mathbf{s}^{(k)})^\top \mathbf{Q} \mathbf{s}^{(k)}} = \frac{(\mathbf{s}^{(k)})^\top \mathbf{s}^{(k)}}{(\mathbf{s}^{(k)})^\top \mathbf{y}^{(k)}} = \alpha_{k+1}^{\text{BB}}$$

★ BB stepsize at  $k+1$  is the SD stepsize at  $k$ .



- ① Derive the DFP formula:

$$\mathbf{H}_{k+1} = \mathbf{H}_k + \frac{\mathbf{s}^{(k)} \mathbf{s}^{(k)\top}}{\mathbf{s}^{(k)\top} \mathbf{y}^{(k)}} - \frac{[\mathbf{H}_k \mathbf{y}^{(k)}] [\mathbf{H}_k \mathbf{y}^{(k)}]^\top}{\mathbf{y}^{(k)\top} \mathbf{H}_k \mathbf{y}^{(k)}}$$

- ② Derive the BFGS formula:

$$\mathbf{B}_{k+1} = \mathbf{B}_k + \frac{\mathbf{y}^{(k)} \mathbf{y}^{(k)\top}}{\mathbf{y}^{(k)\top} \mathbf{s}^{(k)}} - \frac{\mathbf{B}_k \mathbf{s}^{(k)} \mathbf{s}^{(k)\top} \mathbf{B}_k}{\mathbf{s}^{(k)\top} \mathbf{B}_k \mathbf{s}^{(k)}}$$

