No., q	Trial step, $\alpha$		Function value, $f(\alpha)$		
Phase I: Initia	l Bracketing of	Minimum.			
1. $\alpha = 0$	0.000000			3.000000	
2. $q = 0$	$\alpha_0 = \delta =$	$0.500000 \leftarrow \alpha_l$		1.648721	
3. $q = 1$	$\alpha_1 = \sum_{j=0}^1 \delta(1.618)^j = 1.309017$		0.466464		
4. <i>q</i> = 2	$\alpha_2 = \sum_{j=0}^2 \delta_j^2$	$\delta(1.618)^j = 2.6180$	$34 \leftarrow \alpha_n$	5.236610	
Iteration No.	$\alpha_l; [f(\alpha_l)]$	$\boldsymbol{\alpha}_a; [f(\boldsymbol{\alpha}_a)]$	$\boldsymbol{\alpha}_b; [f(\boldsymbol{\alpha}_b)]$	$\alpha_u; [f(\alpha_u)]$	$I = \boldsymbol{\alpha}_u - \boldsymbol{\alpha}_l$
Phase II: Red	ucing Interval o	of Uncertainty			
1	0.500000 [1.648721]↓	1.309017 [0.466 464] レ	1.809017 [0.868376] ↘	2.618034 [5.236610]	2.118034
2	0.500000 [1.648721]	1.000000 ✔ [0.718282]	1.309017 ✔ [0.466464]	1.809017 [0.868376]↓	1.309017
3	1.000000 [0.718282]	1.309017 [0.466 464]	1.500000 [0.481689]	1.809017 [0.868376]	0.809017
_	_	_	_	_	_
_	_	_	_	_	_
16	1.385438 [0.454824]	1.386031 [0.454823]	1.386398 [0.454823]	1.386991 [0.454824]	0.001553
17	1.386031 [0.454823]	1.386398 [0.454823]	1.386624 [0.454823]	1.386991 [0.454823]	0.000960

**TABLE 10.1** Golden Section Search for  $f(\alpha) = 2 - 4\alpha + e^{\alpha}$  of Example 10.3

 $\alpha^* = 0.5(1.386398 + 1.386624) = 1.386511; f(\alpha) = 0.454823.$ 

*Note*: New calculation for each iteration is shown as boldfaced and shaded; the arrows indicate direction of data transfer to the subsequent row/iteration.

#### Inexact Line Search: Armijo's Rule

- Line search termination criterion:  $\nabla f(\mathbf{x}^{(k+1)}) \cdot \mathbf{d}^{(k)} = 0$
- Basic concept
  - step size should not be too large or too small
  - sufficient decrease in the cost function value along the search direction



### Inexact Line Search: Wolfe Conditions

- Sufficient-decrease condition + curvature condition
  - Small values for  $\boldsymbol{\alpha}$

$$f(\alpha) = f\left(\mathbf{x}^{(k)} + \alpha \mathbf{d}^{(k)}\right)$$

$$q(\alpha) = f\left(0\right) + \alpha \left[\rho \underbrace{f'(0)}_{\text{negative}}\right] \rightarrow \begin{cases} \text{sufficient-decrease condition:} f(\alpha) \le q(\alpha) \\ \text{curvature condition:} \\ \frac{\rho \le \beta \le 1}{0 < \rho < 0.5} \Rightarrow f'(\alpha) \ge \beta f'(0) \\ \frac{\beta = 0.9}{\rho = 10^{-4} \sim 10^{-3}} \Rightarrow \left|f'(\alpha)\right| \le \beta \left|f'(0)\right| \Rightarrow \frac{\left|f(\alpha) - f(v)\right|}{\alpha - v} \le \beta \left|f'(0)\right|, \ 0 \le v \le \alpha \\ \text{smaller } \beta, \text{ more accurate } \alpha \rightarrow \begin{cases} \beta = 0.9 : \text{ Newton and quasi-Newton method} \\ \beta = 0.1 : \text{ conjugate gradient method} \end{cases}$$

#### Inexact Line Search: Goldstein Test

- somewhat similar to Armijo's rule
- Newton-type methods, but not quasi-Newton methods



# Multidimensional Unconstrained Optimization

Direct Search Methods	Indirect(Descent) Methods
Random search method	Steepest descent (Cauchy) method
Univariate method	Conjugate gradient method
Pattern search method	<ul> <li>Fletcher-Reeves</li> </ul>
<ul> <li>Powell's method</li> </ul>	– Polak-Rebiere
Simplex method	Newton's method
Simulated Annealing (SA)	Marquardt's method
Genetic Algorithm (GA)	Quasi-Newton methods
	<ul> <li>DFP (Davidon-Fletcher-Powell)</li> </ul>
	– BFGS(Broydon-Fletcher-Goldfarb-Shanno)

### **Properties of Gradient Vector**

- The gradient vector of a function *f* @ the point *x*<sup>\*</sup> is orthogonal (normal) to the tangent plane for the surface *f*=const.
- Gradient represents a direction of maximum rate of increase for the function f at the point  $x^*$ .
- The maximum rate of change of f(x) at any point  $x^*$  is the magnitude of the gradient vector  $\begin{bmatrix} \partial x_1 & \partial x_2 & \partial x_1 \end{bmatrix}$



$$\mathbf{T} = \begin{bmatrix} \frac{\partial x_1}{\partial s} & \frac{\partial x_2}{\partial s} & \cdots & \frac{\partial x_n}{\partial s} \end{bmatrix}$$
$$\frac{df}{ds} = \frac{\partial f}{\partial x_1} \frac{\partial x_1}{\partial s} + \cdots + \frac{\partial f}{\partial x_n} \frac{\partial x_n}{\partial s} = 0$$
$$\frac{df}{dt} = \lim_{\varepsilon \to 0} \frac{f(\mathbf{x} + \varepsilon \mathbf{u}) - f(\mathbf{x})}{\varepsilon}$$
$$= \sum_{i=1}^n u_i \frac{\partial f}{\partial x_i} = (\mathbf{c} \cdot \mathbf{u}) = \mathbf{c}^T \mathbf{u} = \|\mathbf{c}\| \|\mathbf{u}\| \cos \theta$$
$$\max \left| \frac{df}{dt} \right| = \|\mathbf{c}\|$$

Numerical Methods for Unconstrained Optimum Design - 37

# Steepest Descent Method (1)

- Cauchy's Method (1847), first-order method
- Steepest descent direction :  $d = -\nabla f = -\frac{\partial f}{\partial r}$

- satisfy 
$$\nabla f(\mathbf{x}^{(k)}) \cdot \mathbf{d}^{(k)} < 0 \rightarrow - \left\| \nabla f(\mathbf{x}^{(k)}) \right\|^2 < 0$$

- One-dimensional search :  $\frac{df(\mathbf{x}^{(k+1)})}{d\alpha} = \nabla f(\mathbf{x}^{(k+1)}) \cdot \nabla f(\mathbf{x}^{(k)}) = 0$
- ③ simple & robust (convergence guaranteed)
- 🐵 slow rate of convergence
  - Orthogonal path, condition # of the Hessian
- ③ inefficient: no previous information is used

#### **Steepest Descent Algorithm**



# Steepest Descent Method (2)

- Scaling of design variables
  - Accelerate the rate of convergence
  - For a p.d. quadratic function with a unit condition number
    - Converge in only one iteration
  - Unscale the transformed design variables



$$f = x_1^2 + ax_2^2 \rightarrow \nabla^2 f = \begin{bmatrix} 2 & 0 \\ 0 & 2a \end{bmatrix}$$
$$\rightarrow cond(\nabla^2 f) = a$$
$$f = f(\mathbf{x}) \xrightarrow{\mathbf{x} = \mathbf{T}\mathbf{y}} g(\mathbf{y}) = f(\mathbf{T}\mathbf{y})$$
$$\rightarrow \nabla^2 g = \mathbf{T}^T \nabla^2 f \mathbf{T}$$
$$\rightarrow cond(\nabla^2 g) \approx 1$$

# Conjugate Gradient Method (1)

- Conjugate Gradient Direction: Modification of SDM
  - Not orthogonal to each other
  - Cut diagonally through the orthogonal steepest descent directions
  - Orthogonal w.r.t. a symmetric and p.d. A :

$$d^{(i)^{T}}Ad^{(j)} = \begin{cases} a_{jj} & (i = j) \\ 0 & (i \neq j) \end{cases}$$
  
$$d^{(k)} = -\nabla f(\mathbf{x}^{(k)}) + \beta_{k}d^{(k-1)}$$
  
Satisfies descent condition

$$\alpha_{k} \left( \nabla f(\boldsymbol{x}^{(k)}) \cdot \boldsymbol{d}^{(k)} \right) < 0 \rightarrow -\alpha_{k} \nabla f(\boldsymbol{x}^{(k)}) \cdot \nabla f(\boldsymbol{x}^{(k)}) + \alpha_{k} \beta_{k} \underbrace{\nabla f(\boldsymbol{x}^{(k)}) \cdot \boldsymbol{d}^{(k-1)}}_{=0 \text{ (step size condition)}} < 0$$

# Conjugate Gradient Method (2)

- Assumption: exact line search
  - Quadratic function: Fletcher & Reeves(1964)
  - More general function: Polak-Rebiere (1969)
- ③ Simple & Fast Convergent Rate
- *n* iterations for *n* design variables (p.d. quadratic form)

$$\boldsymbol{x}^{(0)} \xrightarrow{\boldsymbol{d}^{(0)}} \boldsymbol{x}^{(1)} \xrightarrow{\boldsymbol{d}^{(1)}} \boldsymbol{x}^{(2)} \cdots \xrightarrow{\boldsymbol{d}^{(n-1)}} \boldsymbol{x}^{(n)}$$

Restart every (n+1) iterations for computational stability

# Conjugate Gradient Method (3)

• Consider the problem of minimizing a quadratic function

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^{T} A \mathbf{x} + \mathbf{c}^{T} \mathbf{x} \to \nabla f(\mathbf{x}^{(k)}) = A \mathbf{x}^{(k)} + \mathbf{c}$$
  

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_{k} d^{(k)} \to f(\alpha) = f(\mathbf{x}^{(k)} + \alpha_{k} d^{(k)})$$
  

$$\frac{\partial f(\alpha)}{\partial \alpha} = d^{(k)^{T}} \nabla f(\mathbf{x}^{(k+1)}) = 0$$
  

$$d^{(k)^{T}} \Big[ A(\mathbf{x}^{(k)} + \alpha_{k} d^{(k)}) + \mathbf{c} \Big] = 0 \to \alpha_{k} = -\frac{d^{(k)^{T}} \nabla f(\mathbf{x}^{(k)})}{d^{(k)^{T}} A d^{(k)}}$$
  

$$d^{(k+1)} = -\nabla f(\mathbf{x}^{(k+1)}) + \beta_{k} d^{(k)}$$
  

$$d^{(k+1)^{T}} A d^{(k)} = -\nabla f(\mathbf{x}^{(k+1)})^{T} A d^{(k)} + \beta_{k} d^{(k)^{T}} A d^{(k)} = 0 \to \beta_{k} = \frac{\nabla f(\mathbf{x}^{(k+1)})^{T} A d^{(k)}}{d^{(k)^{T}} A d^{(k)}}$$

**Optimization Techniques** 

### Polak-Rebiere (1969)

$$\begin{aligned} \mathbf{x}^{(k+1)} &= \mathbf{x}^{(k)} + \alpha_{k} \mathbf{d}^{(k)} \rightarrow \mathbf{d}^{(k)} = \frac{\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}}{\alpha_{k}} \rightarrow \mathbf{A} \mathbf{d}^{(k)} = \frac{\mathbf{A} \mathbf{x}^{(k+1)} - \mathbf{A} \mathbf{x}^{(k)}}{\alpha_{k}} = \frac{\nabla f\left(\mathbf{x}^{(k+1)}\right) - \nabla f\left(\mathbf{x}^{(k)}\right)}{\alpha_{k}} \\ \mathbf{d}^{(k)T} \nabla f\left(\mathbf{x}^{(k+1)}\right) &= 0 \rightarrow \mathbf{d}^{(k-1)T} \nabla f\left(\mathbf{x}^{(k)}\right) = 0 \\ \mathbf{d}^{(k)} &= -\nabla f\left(\mathbf{x}^{(k)}\right) + \beta_{k-1} \mathbf{d}^{(k-1)} \rightarrow \mathbf{d}^{(k-1)} = \frac{\mathbf{d}^{(k)} + \nabla f\left(\mathbf{x}^{(k)}\right)}{\beta_{k-1}} \\ \end{bmatrix} \rightarrow \frac{\left[\mathbf{d}^{(k)} + \nabla f\left(\mathbf{x}^{(k)}\right)\right]^{T} \nabla f\left(\mathbf{x}^{(k)}\right)}{\beta_{k-1}} = 0 \\ \rightarrow \mathbf{d}^{(k)T} \nabla f\left(\mathbf{x}^{(k)}\right) = -\nabla f\left(\mathbf{x}^{(k)}\right)^{T} \nabla f\left(\mathbf{x}^{(k)}\right) \\ \alpha_{k} &= -\frac{\mathbf{d}^{(k)T} \nabla f\left(\mathbf{x}^{(k)}\right)}{\mathbf{d}^{(k)T} \mathbf{A} \mathbf{d}^{(k)}} \rightarrow \alpha_{k} = \frac{\nabla f\left(\mathbf{x}^{(k)}\right)^{T} \nabla f\left(\mathbf{x}^{(k)}\right)}{\mathbf{d}^{(k)T} \mathbf{A} \mathbf{d}^{(k)}} \\ \beta_{k} &= \frac{\nabla f\left(\mathbf{x}^{(k+1)}\right)^{T} \mathbf{A} \mathbf{d}^{(k)}}{\mathbf{d}^{(k)T} \mathbf{A} \mathbf{d}^{(k)}} \rightarrow \beta_{k} = \frac{\nabla f\left(\mathbf{x}^{(k+1)}\right)^{T} \left[\nabla f\left(\mathbf{x}^{(k+1)}\right) - \nabla f\left(\mathbf{x}^{(k)}\right)\right]}{\alpha_{k} \mathbf{d}^{(k)T} \mathbf{A} \mathbf{d}^{(k)}} \\ \rightarrow \overline{\beta_{k}} &= \frac{\nabla f\left(\mathbf{x}^{(k+1)}\right)^{T} \left[\nabla f\left(\mathbf{x}^{(k+1)}\right) - \nabla f\left(\mathbf{x}^{(k)}\right)\right]}{\nabla f\left(\mathbf{x}^{(k)}\right)} \end{bmatrix}$$
for more general objective model

**Optimization Techniques** 

Numerical Methods for Unconstrained Optimum Design - 44

### Fletcher-Reeves (1964)

for a quadratic function 
$$\left[\nabla f\left(\mathbf{x}^{(k+1)}\right) - \nabla f\left(\mathbf{x}^{(k)}\right) = \mathbf{A}\left(\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\right) = \alpha_{k}\mathbf{A}\mathbf{d}^{(k)}\right],$$
  
 $\nabla f\left(\mathbf{x}^{(k+1)}\right)^{T}\nabla f\left(\mathbf{x}^{(k)}\right) = \nabla f\left(\mathbf{x}^{(k+1)}\right)^{T}\left[-\mathbf{d}^{(k)} + \beta_{k-1}\mathbf{d}^{(k-1)}\right] = -\underbrace{\nabla f\left(\mathbf{x}^{(k+1)}\right)^{T}\mathbf{d}^{(k)}}_{=0 \text{ (exact line search)}} + \nabla f\left(\mathbf{x}^{(k+1)}\right)^{T}\beta_{k-1}\mathbf{d}^{(k-1)}$   
 $= \beta_{k-1}\left[\nabla f\left(\mathbf{x}^{(k)}\right)^{T} + \alpha_{k}\mathbf{d}^{(k)T}\mathbf{A}\right]\mathbf{d}^{(k-1)} = \beta_{k-1}\left[\underbrace{\nabla f\left(\mathbf{x}^{(k)}\right)^{T}\mathbf{d}^{(k-1)}}_{=0 \text{ (exact line search)}} + \alpha_{k}\underbrace{\mathbf{d}^{(k)T}\mathbf{A}\mathbf{d}^{(k-1)}}_{=0 \text{ (conjugacy)}}\right] = 0$   
 $\beta_{k} = \frac{\nabla f\left(\mathbf{x}^{(k+1)}\right)^{T}\left[\nabla f\left(\mathbf{x}^{(k+1)}\right) - \nabla f\left(\mathbf{x}^{(k)}\right)\right]}{\nabla f\left(\mathbf{x}^{(k)}\right)} \rightarrow \beta_{k} = \frac{\nabla f\left(\mathbf{x}^{(k+1)}\right)^{T}\nabla f\left(\mathbf{x}^{(k+1)}\right)}{\nabla f\left(\mathbf{x}^{(k)}\right)}$ 

$$d^{(k)} = -\nabla f\left(x^{(k)}\right) + \left(\frac{\left\|\nabla f\left(x^{(k)}\right)\right\|}{\left\|\nabla f\left(x^{(k-1)}\right)\right\|}\right)^2 d^{(k-1)}$$

#### **Conjugate Gradient Algorithm**



#### Example 10.5+6

$$f(x_1, x_2, x_3) = x_1^2 + 2x_2^2 + 2x_3^2 + 2x_1x_2 + 2x_2x_3$$
$$\mathbf{x}^{(0)} = (2, 4, 10), \ \varepsilon = 0.005$$

line search by golden section:  $\delta = 0.05$ ,  $\varepsilon = 0.0001$ 

$f(x_1, x_2, x_3) = x_1^2 + 2x_2^2 + 2x_3^2 + 2x_1x_2 + 2x_2x_3$	
Starting values of design variables	2, 4, 10
Optimum design variables	$8.04787 \hbox{E-} 03, -6.81319 \hbox{E-} 03, 3.42174 \hbox{E-} 03$
Optimum cost function value	2.47347E-05
Norm of gradient of the cost function at optimum	4.97071E-03
Number of iterations	40
Total number of function evaluations	753

TABLE 10.2 Optimum Solution for Example 10.5 with the Steepest–Descent Method

TABLE 10.3 Optimum Solution for Example 10.6 with the Conjugate Gradient Method

$f(x_1, x_2, x_3) = x_1^2 + 2x_2^2 + 2x_3^2 + 2x_1x_2 + 2x_2x_3$	
Starting values of design variables	2, 4, 10
Optimum design variables	1.01E-07, -1.70E-07, 1.04E-09
Optimum cost function value	-4.0E-14
Norm of gradient at optimum	5.20E-07
Number of iterations	4

### Newton's Method

• Optimality criteria for second-order Taylor series

$$f(\mathbf{x} + \Delta \mathbf{x}) = f(\mathbf{x}) + \nabla f^{T}(\mathbf{x})\Delta \mathbf{x} + \frac{1}{2}\Delta \mathbf{x}^{T} \mathbf{H}\Delta \mathbf{x}$$
$$\frac{\partial f}{\partial(\Delta x)} = 0 \Longrightarrow \nabla f(\mathbf{x}) + \mathbf{H}\Delta \mathbf{x} = 0$$
$$\mathbf{d}^{(k)} \equiv \Delta \mathbf{x} = -\mathbf{H}^{-1} \nabla f(\mathbf{x}) \rightarrow \mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \Delta \mathbf{x} \text{ (step length = 1)}$$

- Descent direction?  $\alpha_k (\nabla f(\mathbf{x}^{(k)}) \cdot \mathbf{d}^{(k)}) < 0 \rightarrow -\alpha_k \nabla f(\mathbf{x}^{(k)}) \cdot \mathbf{H}^{-1} \nabla f(\mathbf{x}^{(k)}) < 0$
- Quadratic rate of convergence ③
- Second-order derivative : high cost ⊗
- No guaranteed convergence unless *H* is p.d. ☺
- Memoryless method (no use of previous data) ☺

#### Modified Newton's Algorithm

- Introduce the step length parameter



#### Example 11.8

$$f(\mathbf{x}) = 50(x_2 - x_1^2)^2 + (2 - x_1)^2 \text{ from } \mathbf{x}_0 = (5, -5)$$
  
$$\delta_0 = 0.05$$
  
$$\varepsilon = \begin{cases} 0.00001 \text{ (steepest descent method)} \\ 0.0001 \text{ (modified Newton method)} \end{cases}$$

 TABLE 11.3
 Comparative Evaluation of Three Methods for Example 11.8

$f(x) = 50(x_2 - x_1^2)^2 + (2 - x_1)^2$				
	Steepest-descent	Conjugate gradient	Modified Newton	
<i>x</i> <sub>1</sub>	1.9941E+00	2.0000E+00	2.0000E+00	
<i>x</i> <sub>2</sub>	3.9765E+00	3.9998E+00	3.9999E+00	
f	3.4564E - 05	1.0239E-08	2.5054E - 10	
c	3.3236E-03	1.2860E - 04	9.0357E-04	
Number of function evaluations	138236	65	349	
Number of iterations	9670	22	13	

# Marquardt Modification

- Marquardt (1963)
  - Steepest descent method: far away from the solution
  - Newton method: near the solution point

$$\boldsymbol{d}^{(k)} = -(\boldsymbol{H} + \lambda \boldsymbol{I})^{-1} \nabla f(\boldsymbol{x})$$

- When  $\lambda$  is large,
  - steepest descent direction:  $d^{(k)} = (-1/\lambda) \nabla f(x)$
- As  $\lambda$  is reduced (step size is increased),
  - Newton direction:  $d^{(k)} = -H^{-1}\nabla f(x)$
- If the direction does not reduce the cost function, then  $\lambda$  is increased (step size is reduced)

#### Marquardt's Algorithm (1)



# Marquardt's Algorithm (2)

- Advantages
  - Simplicity
  - Descent property
  - Excellent convergence rate near x<sup>\*</sup>
  - Absence of a line search
- Disadvantages
  - Need to calculate  $H^{(k)}$
  - Solve the linear equation set

$$(\boldsymbol{H} + \lambda^{(k)}\boldsymbol{I})\boldsymbol{d}^{(k)} = -\nabla f(\boldsymbol{x}^{(k)})$$

### **Quasi-Newton Method**

- No learning process
  - Steepest Descent Method: Convergent, but Slow
  - Newton's Method: Fast, but Expensive
- Use of previous information, speed up the convergence !
  - Approximate Hessian (or its inverse) matrix by 1st-order derivatives preserving symmetry and positive definiteness
- Variable Metric Method

$$\boldsymbol{d}^{(k)} = -\boldsymbol{A}^{(k)} \nabla f\left(\boldsymbol{x}^{(k)}\right)$$
$$\boldsymbol{A}^{(k+1)} = \boldsymbol{A}^{(k)} + \boldsymbol{A}_{c}^{(k)} \xrightarrow{\text{as } k \to \infty} \boldsymbol{H}^{-1}$$

# DFP Method (1)

- Davidon (1959)  $\rightarrow$  Fletcher and Powell (1963)
  - Approximate inverse of Hessian matrix
    - For a quadratic objective, directions of the conjugate gradient method
  - One of the most powerful algorithm for general function
  - $\otimes$  need to store the *N*×*N* matrix *A*

$$\begin{aligned} x^{(k+1)} &= x^{(k)} + \alpha_k d^{(k)} \\ d^{(k)} &= -A^{(k)} \nabla f \left( x^{(k)} \right) \\ A^{(k+1)} &= A^{(k)} + \frac{s^{(k)} s^{(k)^T}}{s^{(k)^T} y^{(k)}} - \frac{z^{(k)} z^{(k)^T}}{y^{(k)^T} z^{(k)}} \\ s^{(k)} &= \alpha_k d^{(k)} = x^{(k+1)} - x^{(k)} \\ y^{(k)} &= \nabla f \left( x^{(k+1)} \right) - \nabla f \left( x^{(k)} \right) \\ z^{(k)} &= A^{(k)} y^{(k)} \end{aligned}$$

$$x^{(k+1)} = x^{(k)} + \alpha_k d^{(k)}$$
  

$$\to \Delta x^{(k)} = -\alpha_k A^{(k)} \nabla f(x^{(k)})$$
  

$$c^{(k)} = \nabla f(x^{(k)})$$
  

$$f(x^{(k+1)}) = f(x^{(k)}) + c^{(k)^T} \Delta x^{(k)}$$
  

$$\Delta f(x^{(k)}) = -\alpha^{(k)} c^{(k)^T} A^{(k)} c^{(k)} < 0$$

### DFP Method (2)



# BFGS Method (1)

- Broyden-Fletcher-Goldfarb-Shanno (1981)
  - Direct update the Hessian matrix
  - Most popular and effective

$$\begin{aligned} x^{(k+1)} &= x^{(k)} + \alpha_k d^{(k)} \\ H^{(k)} d^{(k)} &= -\nabla f \left( x^{(k)} \right) \\ H^{(k+1)} &= H^{(k)} + \frac{y^{(k)} y^{(k)^T}}{y^{(k)^T} s^{(k)}} - \frac{H^{(k)} s^{(k)} s^{(k)^T} H^{(k)}}{s^{(k)^T} H^{(k)} s^{(k)}} \\ &\xrightarrow{s^{(k)} = x^{(k+1)} - x^{(k)} = \alpha_k d^{(k)}}_{H^{(k+1)} = -\alpha_k c^{(k)}} H^{(k+1)} = H^{(k)} + \frac{y^{(k)} y^{(k)^T}}{y^{(k)^T} s^{(k)}} + \frac{c^{(k)} c^{(k)^T}}{c^{(k)^T} d^{(k)}} \\ s^{(k)} &= \alpha_k d^{(k)} = x^{(k+1)} - x^{(k)} \\ y^{(k)} &= c^{(k+1)} - c^{(k)} = \nabla f \left( x^{(k+1)} \right) - \nabla f \left( x^{(k)} \right) \end{aligned}$$

### BFGS Method (2)



# Remarks

- DFP and BFGS methods have theoretical properties that guarantee
  - Superlinear (fast) convergence rate
  - Global convergence under certain conditions
- Duals of each other:  $s \leftrightarrow y$ ,  $H \leftrightarrow A$
- Both methods could fail for nonlinear problems,
  - DFP is highly sensitive to inaccuracies in line searches.
  - Both methods can get stuck on a saddle point.
  - Update of Hessian becomes "corrupted" by round-off and other inaccuracies.
- All kinds of "tricks" such as scaling and preconditioning exist to boost the performance of the methods

### **Gradient-Based Methods**

Method	Direction
Steepest Descent	$\boldsymbol{d}^{(k)} = -\nabla f(\boldsymbol{x}^{(k)})$
Conjugate Gradient	$\boldsymbol{d}^{(k)} = -\nabla f(\boldsymbol{x}^{(k)}) + \beta_k \boldsymbol{d}^{(k-1)} \text{ where } \beta_k = \left\  \nabla f(\boldsymbol{x}^{(k)}) \right\ ^2 / \left\  \nabla f(\boldsymbol{x}^{(k-1)}) \right\ ^2$
Newton's	$\boldsymbol{d}^{(k)} = -\boldsymbol{H}^{-1} \nabla f\left(\boldsymbol{x}^{(k)}\right)$
Quasi-Newton	DFP: $d^{(k)} = -A\nabla f(\mathbf{x}^{(k)})$ where $A^{(k+1)} = A^{(k)} + \frac{s^{(k)}s^{(k)^{T}}}{s^{(k)^{T}}y^{(k)}} - \frac{z^{(k)}z^{(k)^{T}}}{y^{(k)^{T}}z^{(k)}}$ BFGS: $H^{(k)}d^{(k)} = -\nabla f(\mathbf{x}^{(k)})$ where $H^{(k+1)} = H^{(k)} + \frac{y^{(k)}y^{(k)^{T}}}{y^{(k)^{T}}s^{(k)}} + \frac{c^{(k)}c^{(k)^{T}}}{c^{(k)^{T}}d^{(k)}}$

# Rate/Order of Convergence

• Algorithm has p order of convergence if

$$0 \le \beta = \lim_{k \to \infty} \frac{\left\| x^{(k+1)} - x^* \right\|}{\left\| x^{(k)} - x^* \right\|^p} < \infty$$

 $\beta$ : convergence ratio (asymptotic error constant)

- Linear convergence rate: p = 1,  $\beta \le 1$ 
  - Steepest Descent Method
- Superlinear convergence: p = 1,  $\beta = 0$ 
  - Quasi-Newton Method
- Quadratic convergence: p = 2
  - Newton's Method