

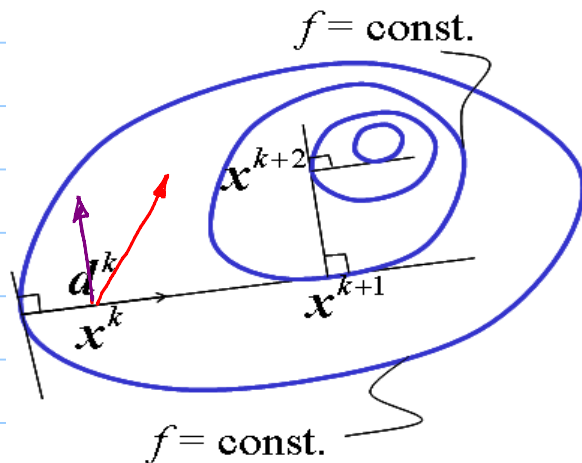
# 3-3 : Conjugate Gradient Method

노트 제목

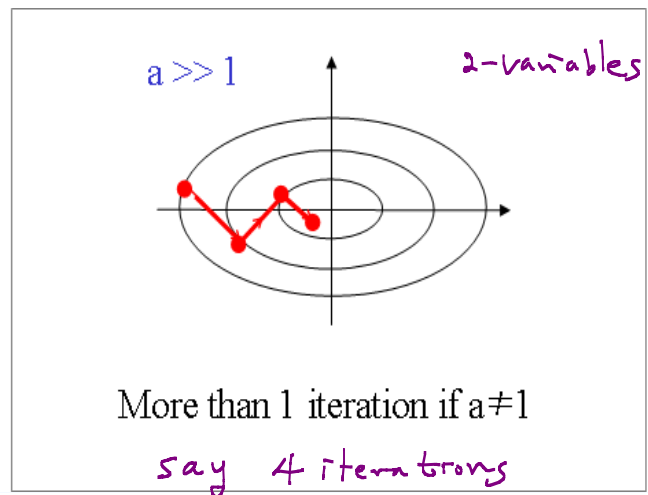
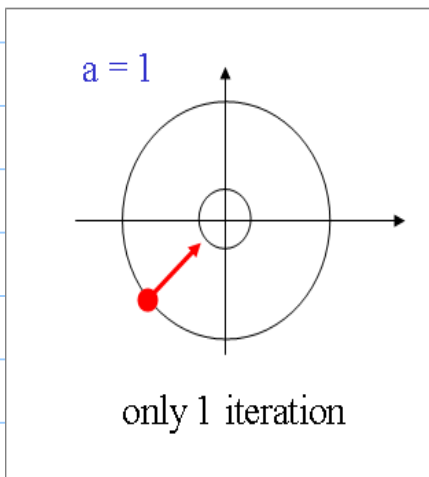
①

< Motivation >

i) Recall the convergence pattern of the steepest descent method



ii) Consider  $f(x_1, x_2) = x_1^2 + ax_2^2$  ← Quadratic ftn of 2 variables



⇒ We want a method yielding a min of  $f(x)$  within "n" iterations if  $f(x)$  is a quadratic function of "n" variables

(2)

⇒ "Conjugate Gradient Method"

(a very powerful method that also works for general non-quadratic functions)

Will study

1) What are conjugate directions?

i.e., definition of conjugate direction

2) Proof to show  $n$ -iteration convergence of  $Q(x_1, \dots, x_n)$  with conjugate search direction.

3) How to find conjugate direction?  
(approximate for non-quadratic fns)

(1) Definition of Conjugate Directions

Given  $\underline{A}$ :  $n \times n$  positive-definite symmetric matrix

$\underline{d}^i$  : conjugate direction

if  $(\underline{d}^i)^T \underline{A} \underline{d}^j = 0$  for  $i \neq j$

Examples of conjugate directions

Let  $A$ :  $3 \times 3$  sym matrix

$e_i$ : eigenvector of  $A$  ( $e_i \cdot A e_j \equiv \delta_{ij}$ )

Then

$$\left\{ \begin{array}{l} (e_1, e_2, e_3) : \text{conjugate direction} \\ (e_1 + e_2, e_1 - e_2, e_3) : \text{''} \\ \vdots \end{array} \right.$$

\* A special case of conjugate directions = eigenvectors

(2) < Proof of n-iteration Convergence >

Min of a quadratic ftn of  $n$  variables is found in  $n$  iterations or less.

Consider  $Q(x) = \frac{1}{2} x^T A x + B^T x + C$

$\uparrow$   
 a quadratic ftn

$\uparrow$   
 positive-definite symmetric matrix

$x = \{x_1, \dots, x_n\}^T$

(4)

proofLet  $\underline{x}^* = \text{min point}$ 

then  $\underline{\nabla} \varphi(\underline{x}^*) = \underline{A} \underline{x}^* + \underline{B} = 0 \quad (1)$

$$\begin{aligned}
 \left( \frac{\partial \varphi}{\partial x_k} = \frac{\partial}{\partial x_k} \left\{ \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n x_i A_{ij} x_j + \sum_{i=1}^n B_i x_i + C \right\} \right. \\
 &= \frac{1}{2} \left[ \sum_{i=1}^n \sum_{j=1}^n (\delta_{ik} A_{ij} x_j + x_i A_{ij} \delta_{jk}) \right. \\
 &\quad \left. + \sum_{i=1}^n B_i \delta_{ik} \right] \\
 &= \frac{1}{2} \left[ \sum_{j=1}^n A_{kj} x_j + \sum_{i=1}^n x_i A_{ik} \right] + B_k \\
 &= \frac{1}{2} \left[ \sum_{i=1}^n A_{ki} x_i + \sum_{i=1}^n \underbrace{A_{ik}}_{A_{ki}} x_i \right] + B_k \\
 &= \sum_{i=1}^n A_{ki} x_i + B_k
 \end{aligned}$$

Expand  $\underline{x}^*$  in terms of orthogonal  $\underline{d}^i \{i=0, \dots, n-1\}$   
Such that wrt  $\underline{A}$ 

$$\underline{x}^* = \underline{x}^0 + \sum_{i=0}^{n-1} \beta^i \underline{d}^i \quad (2)$$

↑ used to represent  
an arbitrary starting point

⑤

To determine  $\beta^i$ , (2)  $\rightarrow$  (1)

$$\underline{\beta} + \underline{A} \underline{x}^0 + \sum_{i=0}^{n-1} \beta^i \underline{A} \underline{d}^i = 0 \quad (3)$$

$(\underline{d}^j)^T$  (3):

$$\underline{d}^{jT} (\underline{\beta} + \underline{A} \underline{x}^0) + \sum_{i=0}^{n-1} \beta^i \underbrace{\underline{d}^{jT} \underline{A} \underline{d}^i}_{=0 \text{ if } i \neq j} = 0 \quad (4)$$

$\Downarrow$   
 $\beta^j \underline{d}^{jT} \underline{A} \underline{d}^j$

Thus

$$\beta^j = \frac{-(\underline{\beta} + \underline{A} \underline{x}^0)^T \underline{d}^j}{(\underline{d}^j)^T \underline{A} \underline{d}^j} \quad (5)$$

$j = 0, 1, 2, \dots, n-1$

$\Rightarrow$  Message: For "any" starting point  $\underline{x}_0$ , the min  $\underline{x}^*$  of  $Q(\underline{x})$  is expressed in terms of [redacted]  $(\underline{d}^i)$ ; expansion coefficients  $\beta_j$  are computed by (5).

②

Let us now prove that min of quadratic function  $Q(x)$  can be found with  $n$  iterations if we update  $\underline{x}^i$  using  $\underline{d}^i$  as

$$\underline{x}^{i+1} = \underline{x}^i + \lambda^{i*} \underline{d}^i \quad (6)$$

( $i=0, 1, \dots, n-1$ )

↑ conjugate direction

where  $\lambda^{i*}$  is determined by 1-D search by  $\min Q(\underline{x}^i + \lambda^{i*} \underline{d}^i)$

$$\Leftrightarrow \boxed{0 = \frac{dQ}{d\lambda^i} \mid \lambda^i = \lambda^{i*}} \quad (7)$$

\* If we can prove that

$$\lambda^{i*} = \beta^i, \quad (i=0, \dots, n-1)$$

then

$$\underline{x}^1 = \underline{x}^0 + \lambda^{0*} \underline{d}^0 = \underline{x}^0 + \beta^0 \underline{d}^0$$

$$\underline{x}^2 = \underline{x}^1 + \lambda^{1*} \underline{d}^1 = \underline{x}^0 + \beta^0 \underline{d}^0 + \beta^1 \underline{d}^1$$

⋮

(7)

$$\underline{x}^n = \underline{x}^0 + \sum_{i=0}^{n-1} \beta^i \underline{d}^i$$

i.e., updating by (6) and (7) yields solutions within  $n$  iterations.

▣ Let us prove  $\lambda^{i*} = \beta^i$

Use Eq. (7)

$$0 = \frac{dQ(\underline{x}^i + \lambda^i \underline{d}^i)}{d\lambda^i} \Big|_{\lambda^i = \lambda^{i*}}$$

$$= \frac{\partial}{\partial \lambda^i} \left[ \frac{1}{2} (\underline{x}^i + \lambda^i \underline{d}^i)^T \underline{A} (\underline{x}^i + \lambda^i \underline{d}^i) + \underline{B}^T (\underline{x}^i + \lambda^i \underline{d}^i) + C \right]_{\lambda^i = \lambda^{i*}}$$

$$= \underline{d}^{iT} \underline{A} (\underline{x}^i + \lambda^{i*} \underline{d}^i) + \underline{B}^T \underline{d}^i$$

$$= \lambda^{i*} (\underline{d}^{iT} \underline{A} \underline{d}^i) + \underbrace{\underline{d}^{iT} \underline{A} \underline{x}^i + \underline{B}^T \underline{d}^i}_{\text{Scalar}}$$

$$= \underline{x}^{iT} \underline{A}^T \underline{d}^i$$

$$\underline{x}^{iT} \underline{A} \underline{d}^i$$

(8)

$$\therefore \lambda^{i*} = - \frac{(\underline{b}^T + \underline{x}^{iT} \underline{A}) \underline{d}^i}{\underline{d}^{iT} \underline{A} \underline{d}^i} \quad (8)$$

To simplify (8), note that

$$\begin{aligned} \textcircled{1} \quad \underline{x}^i &= \underline{x}^{i-1} + \lambda^{(i-1)*} \underline{d}^{i-1} \\ &= \underline{x}^{i-2} + \lambda^{(i-2)*} \underline{d}^{i-2} + \lambda^{(i-1)*} \underline{d}^{i-1} \\ &= \underline{x}^0 + \sum_{k=0}^{i-1} \lambda^{k*} \underline{d}^k \end{aligned}$$

$$\begin{aligned} \textcircled{2} \quad \underline{x}^{iT} \underline{A} \underline{d}^i &= (\underline{x}^0 + \sum_{k=0}^{i-1} \lambda^{k*} \underline{d}^k)^T \underline{A} \underline{d}^i \\ &= \underline{x}^{0T} \underline{A} \underline{d}^i + \sum_{k=0}^{i-1} \lambda^{k*} \underbrace{\underline{d}^{kT} \underline{A} \underline{d}^i}_{=0} \\ &= \underline{x}^{0T} \underline{A} \underline{d}^i \quad \text{Because } \underline{d}^k, \text{ conjugate} \\ &= \underline{x}^{0T} \underline{A}^T \underline{d}^i \quad \text{directions} \\ &= (\underline{A} \underline{x}^0)^T \underline{d}^i \quad \left. \begin{array}{l} \\ \\ \end{array} \right\} i \neq j \end{aligned}$$

$$\therefore \lambda^{i*} = - \frac{(\underline{b} + \underline{A} \underline{x}^0)^T \underline{d}^i}{\underline{d}^{iT} \underline{A} \underline{d}^i} \equiv \beta^i$$



(3) < How to find  $\underline{d}$ ? >

Result:

□ for quadratic function  $Q(\underline{x}) = \frac{1}{2} \underline{x}^T \underline{A} \underline{x} + \underline{B}^T \underline{x} + C$

$$\bullet \underline{d}^0 = -\underline{g}^0 \quad (\underline{g} = \nabla Q)$$

$$\bullet \underline{d}^{k+1} = -\underline{g}^{k+1} + \beta^k \underline{d}^k \quad (k \geq 1)$$

$$\text{where } \beta^k = \frac{(\underline{g}^{k+1})^T \underline{A} \underline{d}^k}{\underline{d}^{kT} \underline{A} \underline{d}^k}$$

□ for non-quadratic function  $f(\underline{x})$

replace  $\beta^k$  as

$$\beta^k = \frac{\underline{g}^{k+1T} \underline{g}^{k+1}}{\underline{g}^{kT} \underline{g}^k}$$

(Key point: no second-derivative information such as Hessian ( $\approx \underline{A}$ ) is needed  
→ due to Fletcher-Reeve)

overall approach

- i) derive  $\underline{d}$  for  $Q(\underline{x})$
- ii) extend the result for non-quadratic functions

## &lt; Derivation &gt;

(A) For  $Q(\underline{x}) = \frac{1}{2} \underline{x}^T \underline{A} \underline{x} + \underline{B}^T \underline{x} + C$  first

- $\underline{d}_0$  can be any non-zero vector to be a conjugate vector at the beginning

• thus, choose  $\underline{d}_0 = -\underline{g}_0$

where

$$\begin{aligned} \underline{g}^k &= \underline{\nabla} Q(\underline{x}^k) \\ &= \underline{A} \underline{x}^k + \underline{B} \end{aligned} \quad (a)$$

- Search Algorithm

$$\underline{x}^{k+1} = \underline{x}^k + \alpha^k \underline{d}^k \quad (b)$$

↑  
found by 1-D search

- Conjugate direction  $\underline{d}^k$

$$\left\{ \begin{array}{l} \underline{d}^k = -\underset{\substack{\uparrow \\ \text{known}}}{\underline{g}^k} + \beta^{k+1} \underset{\substack{\uparrow \\ \text{unknown}}}{\underline{d}^{k+1}} \quad (c1) \\ \text{or} \\ \underline{d}^{k+1} = -\underline{g}^{k+1} + \beta^k \underline{d}^k \quad (c2) \end{array} \right.$$

①

- To determine  $\beta^k$ , use the conjugate condition.

$$\underline{d}^{kT} \underline{A} \underline{d}^{k+1} = 0$$

$$[-\underline{g}^k + \beta^{k+1} \underline{d}^{k+1}]^T \underline{A} \underline{d}^{k+1} = 0$$

$$\rightarrow \beta^{k+1} = \frac{(\underline{g}^k)^T \underline{A} \underline{d}^{k+1}}{(\underline{d}^{k+1})^T \underline{A} \underline{d}^{k+1}}$$

or

$$\beta^k = \frac{(\underline{g}^k)^T \underline{A} \underline{d}^k}{(\underline{d}^k)^T \underline{A} \underline{d}^k} \quad (d)$$

for non-quadratic fcn,  
 $\underline{A}$  can be calculated logically  
as  $\underline{H}(\underline{x}^k)$ , but calculation  
of  $\underline{H}(\underline{x}^k)$  should be avoided  
for computational efficiency

(12)

Extension to Non-Quadratic functions(avoid explicit calculation of  $\underline{A}$ )

Recall

$$\underline{A} \underline{x}^{k+1} = \underline{A} (\underline{x}^k + \alpha^k \underline{d}^k) \quad \leftarrow \text{appearing in (d)}$$

USE (b)

$$\Rightarrow \underline{A} \underline{d}^k = \frac{1}{\alpha^k} (\underline{A} \underline{x}^{k+1} - \underline{A} \underline{x}^k)$$

$$= \frac{1}{\alpha^k} \left[ (\underline{A} \underline{x}^{k+1} + \underline{B}) - (\underline{A} \underline{x}^k + \underline{B}) \right]$$

$$\stackrel{\text{USE (a)}}{=} \frac{1}{\alpha^k} (\underline{g}^{k+1} - \underline{g}^k) \quad \text{(e)}$$

Expressed only in terms of  
gradients(e)  $\rightarrow$  (d) :

$$\beta^k = \frac{(\underline{g}^{k+1})^T (\underline{g}^{k+1} - \underline{g}^k)}{(\underline{d}^k)^T (\underline{g}^{k+1} - \underline{g}^k)} \quad \text{(f)}$$

To simplify (f), use

$$\begin{cases} \textcircled{1} (\underline{g}^{k+1})^T \underline{d}^k = 0 \\ \textcircled{2} (\underline{d}^k)^T \underline{g}^k = -(\underline{g}^k)^T \underline{g}^k \\ \textcircled{3} (\underline{g}^{k+1})^T \underline{g}^k = 0 \end{cases}$$

Then

$$\beta^k \stackrel{(f)}{=} \frac{(\underline{g}^{k+1})^T \underline{g}^{k+1} - (\underline{g}^{k+1})^T \underline{g}^k}{(\underline{d}^k)^T \underline{g}^{k+1} - (\underline{d}^k)^T \underline{g}^k} \quad (13)$$

by ①
by ②

$$\beta^k = \frac{(\underline{g}^{k+1})^T \underline{g}^k}{(\underline{g}^k)^T \underline{g}^k}$$

## proofs

proof of ①:  $(\underline{g}^{k+1})^T \underline{d}^k = 0$

Recall:  $f(\underline{x}) = f(\underbrace{\underline{x}^k}_{\text{known}} + \alpha \underbrace{\underline{d}^k}_{\text{known}})$

→ becomes min at  $\alpha = \alpha^k$  such that

$$\begin{aligned} 0 &= \left. \frac{df}{d\alpha} \right|_{\alpha=\alpha^k} = \left[ (\nabla f)^T \underline{d}^k \right]_{\alpha=\alpha^k} \\ &= \left[ \nabla f(\underbrace{\underline{x}^k + \alpha^k \underline{d}^k}_{\underline{x}^{k+1}}) \right]^T \underline{d}^k \\ &= (\underline{g}^{k+1})^T \underline{d}^k \end{aligned}$$

(14)

proof of ② :  $(\underline{d}^k)^T \underline{g}^k = -(\underline{g}^k)^T \underline{g}^k$

$$\begin{aligned} (\underline{d}^k)^T \underline{g}^k &\stackrel{(c)}{=} [-\underline{g}^k + \beta^{k+1} \underline{d}^{k-1}]^T \underline{g}^k \\ &= -(\underline{g}^k)^T \underline{g}^k + \beta^{k+1} \underbrace{(\underline{d}^{k-1})^T \underline{g}^k}_{\stackrel{||}{=} (\underline{g}^k)^T \underline{d}^{k-1} \equiv 0 \text{ by } \textcircled{D}} \\ &= -(\underline{g}^k)^T \underline{g}^k \end{aligned}$$

proof of ③ :  $(\underline{g}^{k+1})^T \underline{g}^k = 0$

$$\begin{aligned} (\underline{g}^{k+1})^T \underline{g}^k &\stackrel{(c)}{=} (\underline{g}^{k+1})^T [-\underline{d}^k + \beta_{k+1} \underline{d}^{k-1}] \\ &= -\underbrace{(\underline{g}^{k+1})^T \underline{d}^k}_{\equiv 0 \text{ by } \textcircled{D}} + \beta_{k+1} (\underline{g}^{k+1})^T \underline{d}^{k-1} \\ &= \beta_{k+1} (\underline{g}^{k+1})^T \underline{d}^{k-1} \quad (g) \end{aligned}$$

Now rewrite (g) as

$$\underline{g}^{k+1} = \underline{g}^k + \alpha^k \underline{A} \underline{d}^k \quad (h)$$

Then

$$\begin{aligned} \underline{g}^{k+1} &\stackrel{(g)}{=} \beta_{k+1} \underbrace{[\underline{g}^k + \alpha_k \underline{A} \underline{d}^k]}_{(h)}^T \underline{d}^{k+1} \\ &= \beta_{k+1} \underbrace{(\underline{g}^k)^T \underline{d}^{k+1}}_{\text{by } \textcircled{D}} + \beta_{k+1} \alpha_k \underbrace{(\underline{d}^k)^T \underline{A} \underline{d}^{k+1}}_{\text{by conjugate cond.}} = 0 \end{aligned}$$

15

## < Fletcher-Reeve Conjugate Gradient Method >

- Start with any  $\underline{x}_0$

- Search direction

if  $k = 0$

$$\underline{d}^k = -\underline{g}^k = -\nabla f(\underline{x}^k)$$

else

$$\underline{d}^k = -\underline{g}^k + \beta^{k-1} \underline{d}^{k-1}$$

$$\beta^{k-1} = \frac{(\underline{g}^k)^T \underline{g}^k}{(\underline{g}^{k-1})^T \underline{g}^{k-1}} = \frac{\|\nabla f(\underline{x}^k)\|^2}{\|\nabla f(\underline{x}^{k-1})\|^2}$$

- Repeat until convergence

Remark: ① this method applies to non-quadratic function minimization

②  $n$ -iteration convergence is valid only for quadratic ftn

③ may need to reset  $\underline{d}^k$  for every  $n$  step

## Example

Consider  $f = x_1^2 + 4x_2^2$ ,  $\mathbf{x}_0 = (1,1)^T$

- the steepest descent iteration

$$\mathbf{d}_0 = -\nabla f(\mathbf{x}_0) = -(2, 8)^T$$

$$f(\alpha) = f(\mathbf{x}_0 + \alpha \mathbf{d}_0) \Rightarrow \alpha_0 = 0.1308$$

$$\therefore \mathbf{x}_1 = \mathbf{x}_0 + \alpha_0 \mathbf{d}_0 = (0.7385, -0.0462)^T$$

- the conjugate gradient iteration

$$\beta_0 = \frac{\|\nabla f(\mathbf{x}_1)\|^2}{\|\nabla f(\mathbf{x}_0)\|^2} = 0.0341$$

$$\mathbf{d}_1 = -\nabla f^T(\mathbf{x}_1) + \beta_0 \mathbf{d}_0 = \begin{pmatrix} -1.5451 \\ 0.0966 \end{pmatrix}$$

$$f(\alpha) = f(\mathbf{x}_1 + \alpha \mathbf{d}_1) \Rightarrow \alpha_1 = 0.4779$$

$$\therefore \mathbf{x}_2 = \mathbf{x}_1 + \alpha_1 \mathbf{d}_1 \cong (0.0, 0.0)^T$$

