# Lecture 3-2  The Steepest Descent Method



▣ Consider the update rule
$$\underset{\sim}{x}^{k+1} = \underset{\sim}{x}^{k} + \alpha \, \underset{\sim}{d}^{k}$$
$$(\text{with } \alpha > 0)$$

to reduce the function value

i.e; $\quad f(\underset{\sim}{x}^{k+1}) < f(\underset{\sim}{x}^{k})$

(repeat until converged)

Question?  For what $\underset{\sim}{d}$, $\quad f(\underset{\sim}{x}^{k+1}) < f(\underset{\sim}{x}^{k})$ valid?

$$f(\underset{\sim}{x}^{k+1}) = f(\underset{\sim}{x}^{k} + \alpha \, \underset{\sim}{d}) \qquad (\alpha > 0)$$

$$\underset{\substack{\text{Taylor}\\ \text{Expansion}}}{=} f(\underset{\sim}{x}^{k}) + \alpha \, \nabla f^{T}(\underset{\sim}{x}^{k}) \, \underset{\sim}{d}$$

$$+ O(\alpha^{2}) \qquad\qquad --- (1)$$

To satisfy
$$f(x^{k+1}) < f(x^{k}) \qquad (\alpha > 0) \qquad -- (2)$$

$\underset{\sim}{d}$ must satisfy

$$\boxed{\nabla f^{T}(\underset{\sim}{x}_{k}) \cdot \underset{\sim}{d} < 0}$$

(3)

**Consequence of (3)**

$$\nabla f : \text{Function-Increasing direction} \quad \text{---(A)}$$

**Claim** (will be proven)

$\underline{\nabla} f$ is orthogonal to the contour (B)
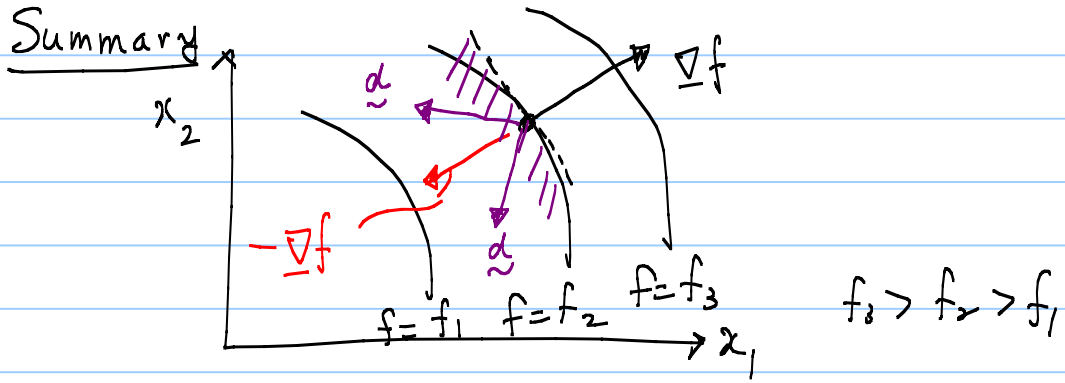
of $f = \text{const}$

$$\underline{\nabla} f^T \underline{t} = 0$$

$$(\underline{\nabla} f \cdot \underline{t} = 0)$$

$f = f_1 < f = f_2 < f = f_3$

($\underline{t}$ : tangent vector)

$(A, B) \Rightarrow \underset{\sim}{\nabla} f :$ the direction of fastest function increase $\rightarrow$ Steepest ascent direction

$\therefore -\underset{\sim}{\nabla} f :$ Steepest descencent direction

Summary



$f_3 > f_2 > f_1$

i) $f$ decreases for any $\underset{\sim}{d}$ lying in the shaded region

ii) steepest descent direction $\underset{\sim}{d}$

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

$$\underset{\text{Normalize}}{=} -\frac{\underset{\sim}{\nabla} f(\underline{x}^k)}{\|\underset{\sim}{\nabla} f(\underline{x}^k)\|}$$

< Proof of Claim B >



Along $f = $ CONST

$0 = df/ds$  ($s$: arclength along $C$)

$$= \left[\frac{\partial f}{\partial x}\frac{dx}{ds} + \frac{\partial f}{\partial y}\frac{dy}{ds}\right]_{\text{along } C'}$$

$$= \left(\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}\right) \cdot \left(\frac{dx}{ds}, \frac{dy}{ds}\right)$$

$$= \underline{\nabla f}^T \underset{\sim}{t}$$

For the steepest descent Method

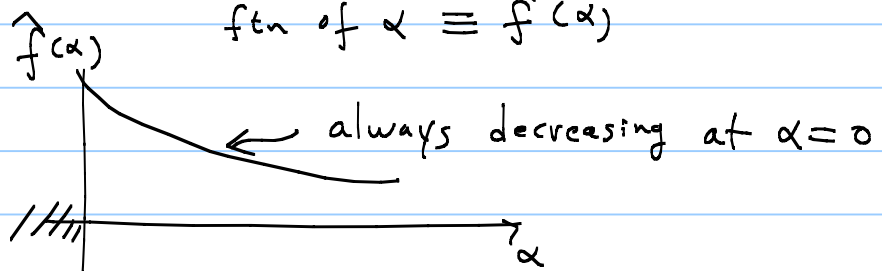i) Search direction $\underset{\sim}{d}_k = -\nabla f(\underline{x}_k) / \|\nabla f(\underline{x}_k)\|$

ii) Step $\alpha > 0$ ?

$$\underline{x}^{k+1} = \underline{x}^k + \alpha \underline{d}^k$$

$$\Rightarrow \underset{\sim}{x}^k = \underline{x}^{k+1}(\alpha) \; ; \; \text{1-D problem}$$

$$\therefore \underset{\alpha}{\text{Min}} \; \underbrace{f(\underline{x}^k + \alpha \underset{\sim}{d}^k)}_{\text{ftn of } \alpha \equiv \hat{f}(\alpha)}$$

Remark $\hat{f}(\alpha)$



← always decreasing at $\alpha = 0$

Check $\left. \dfrac{d\hat{f}}{d\alpha} \right|_{\alpha=0} = \left. \dfrac{d}{d\alpha} f(\underline{x}^k + \alpha \underline{d}^k) \right|_{\alpha=0}$

$\underline{\underline{\triangleq}} \underset{\sim}{x}$

$\left( x_i = x_i^k + \alpha d_i^k \right)$

$$= \sum_{i=1}^{n} \left[ \frac{\partial f}{\partial x_i} \frac{\partial}{\partial \alpha} x_i \right]_{\alpha=0}$$

$$= \sum_{i=1}^{n} \left. \frac{\partial f}{\partial x_i} \right|_{\underset{\sim}{x}^k} d_i^k = \nabla f^T(\underline{x}^k) \, \underline{d}^k < 0$$

Recall how $\underline{d}$ was chosen in the steepest descent method

**5** <u>Stopping criteria</u>

   i) check the necessary condition for
      optimality     $\| \nabla f(\underline{x}^k) \| \leq \epsilon_G$
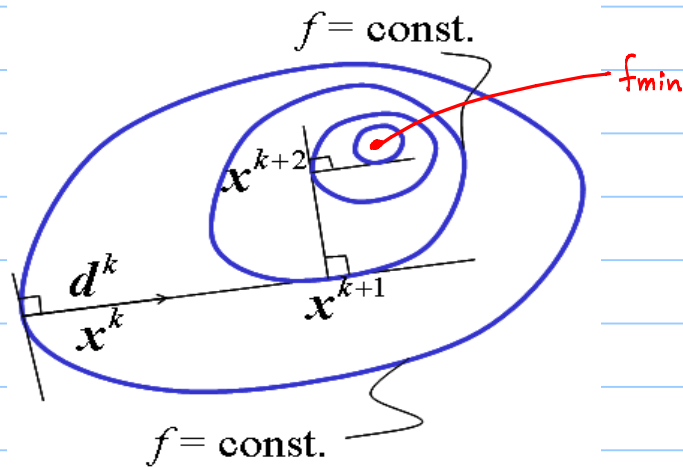
$\qquad\qquad\qquad$ ( usually $O(10^{-6})$ )

   ii) Check the sucessive reduction in $f$

   $| f(\underline{x}^{k+1}) - f(\underline{x}^k) | \leq \epsilon_A + \epsilon_R | f(\underline{x}^k) |$

$\qquad$ usually $\epsilon_A, \epsilon_R = O(10^{-6})$

   (Check at least two sucessive iterations before
   the stopping the process )


**6** Convergence property = ?



$f$ = const.

fmin

$x^{k+2}$

$d^k$

$x^k$

$x^{k+1}$

$f$ = const.

① $\boxed{\underset{\sim}{d}^{K+1} \perp \underset{\sim}{d}^{K} \quad \left( \Leftrightarrow \left(\underset{\sim}{d}^{K+1}\right)^{T} \underset{\sim}{d}^{K} = 0 \right)}$

$\rightarrow$ Every Search direction is orthogonal to the previous step

proof: Consider $f(\underset{\sim}{x}^{k+1}) = f(\underset{\sim}{x}^{k} + \alpha^{k} \underset{\sim}{d}^{k})$

with $\dfrac{d}{d\alpha} f(\underset{\sim}{x}^{k} + \alpha \underset{\sim}{d}^{k}) \Big|_{\alpha = \alpha_{K}} = 0$     (a)

(by 1-D Search)

Due to (a)

$0 = \dfrac{d}{d\alpha} f(\underbrace{\underset{\sim}{x}^{k} + \alpha \underset{\sim}{d}^{k}}_{\underset{\sim}{x}}) \Big|_{\alpha = \alpha_{K}}$

$\left( x_{i} = x_{i}^{k} + \alpha d_{i}^{k} \right) = \left[ \displaystyle\sum_{i=1}^{n} \dfrac{\partial f}{\partial x_{i}} \dfrac{\partial x_{i}}{\partial \alpha} \right]_{\alpha = \alpha_{K}}$

$= \displaystyle\sum_{i=1}^{n} \dfrac{\partial f}{\partial x_{i}} \Big|_{x^{k+1}} d_{i}^{k} = \nabla f^{T}(\underset{\sim}{x}^{k+1}) \underset{\sim}{d}^{k}$

Because $\underset{\sim}{d}^{k+1} = -\nabla f(\underset{\sim}{x}^{k+1})$,

$\boxed{(d^{K+1})^{T} \underset{\sim}{d}^{k} = 0} \quad \Longleftarrow$ However, this relation is valid when the 1-D line search is exact.

② What controls the convergence rate of the steepest descent method?

OBSERVATION: $f(x_1, x_2) = x_1^2 + a x_2^2$



a = 1

only 1 iteration

a >> 1

More than 1 iteration if a≠1

Condition number of Hessian Matrix $\underset{\sim}{H}$
(i.e, $|\lambda_{max}| \lambda_{min}|$ of $\underline{H}$) affects
the convergence property

Consider : $f(x_1, x_2) = x_1^2 + a x_2^2$

■ Check the Hessian of $f$ :

- $H = \left[ \dfrac{\partial^2 f}{\partial x_i \partial x_j} \right] = \begin{bmatrix} 2 & 0 \\ 0 & 2a \end{bmatrix}$
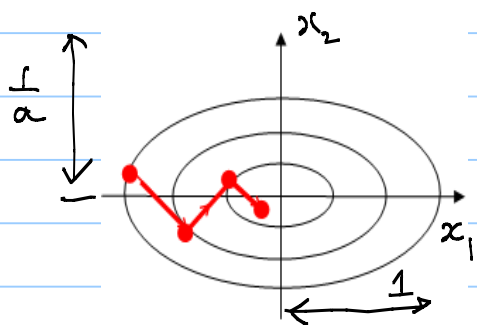
- Eigenvalue of $H$

$$HZ - \lambda IZ = 0 \qquad \begin{bmatrix} 2-\lambda & 0 \\ 0 & 2a-\lambda \end{bmatrix} \begin{Bmatrix} z_1 \\ z_2 \end{Bmatrix} = \begin{Bmatrix} 0 \\ 0 \end{Bmatrix}$$

$$det. = 0 \qquad \rightarrow \qquad \lambda = 2, \ \lambda = 2a$$

$$if \ a > 1, \qquad Condition \ Number = \frac{\lambda_{max}}{\lambda_{min}} = \frac{2a}{2} = a$$

‹Trick to Improve Convergence ?›

$$f = x_1^2 + a\, x_2^2 \qquad\qquad f = y_1^2 + y_2^2$$



Scaling!!

$$y_1 = x_1$$
$$y_2 = \sqrt{a}\, x_2$$

$$\begin{Bmatrix} x_1 \\ x_2 \end{Bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & \frac{1}{\sqrt{a}} \end{bmatrix} \begin{Bmatrix} y_1 \\ y_2 \end{Bmatrix} \overset{\triangle}{=} \underset{\sim}{T} \begin{Bmatrix} y_1 \\ y_2 \end{Bmatrix}$$

$$\boxed{\underset{\sim}{x}^k = \underset{\sim}{D}^k \underset{\sim}{y}^k} \qquad \leftarrow \text{Transform } \underset{\sim}{x} \text{ to } \underset{\sim}{y}$$

$$\underset{\sim}{H}_x = \begin{bmatrix} \frac{\partial^2 f}{\partial x_i \partial x_j} \end{bmatrix} \quad , \quad \underset{\sim}{H}_y = \begin{bmatrix} \frac{\partial^2 f}{\partial y_i \partial y_j} \end{bmatrix}$$

$$\frac{\partial^2 f}{\partial y_i \partial y_j} = D_i \frac{\partial^2 f}{\partial x_i \partial x_j} D_j$$

$$\underset{\sim}{H}_y = \underset{\sim}{D}^T \underset{\sim}{H}_x \underset{\sim}{D}$$

$$= \begin{bmatrix} 1 & 0 \\ 0 & \frac{1}{\sqrt{a}} \end{bmatrix} \begin{bmatrix} 2 & 0 \\ 0 & 2a \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & \frac{1}{\sqrt{a}} \end{bmatrix} = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$

$$\text{Cond}(H_y) = 1.$$

More General Approch:

$$\underset{\sim}{x} = \underset{\sim}{T} \underset{\sim}{y}$$

several approches to Choose $\underset{\sim}{T}$ available

can be non-diagonal matrix though Diagonal matrix is common

▣ Application of $\underset{\sim}{x} = \underset{\sim}{T} \underset{\sim}{y}$

$$\boxed{\text{Solve} \quad A\underset{\sim}{x} = \underset{\sim}{b} \quad \text{for "very" large } n}$$

① Convert it as a minimization problem:

$$\min_{\underset{\sim}{x}} Q = \frac{1}{2} \underset{\sim}{x}^T \underset{\sim}{A} \underset{\sim}{x} - \underset{\sim}{b}^T \underset{\sim}{x} \qquad (A_{ij} = A_{ji}^T)$$

$\left(\rule{0pt}{60pt}\right.$

• $Q = \frac{1}{2} \sum_i \sum_j A_{ij} x_i x_j - \sum_i b_i x_j$

• NC for $Q$ to be min

$\dfrac{\partial Q}{\partial x_k} = 0 : \dfrac{1}{2}\left[ \sum_i \sum_j A_{ij} \underbrace{\left(\dfrac{\partial x_i}{\partial x_k}\right)}_{\delta_{ik}} x_j \right.$

$\left. + \sum_i \sum_j A_{ij} x_i \underbrace{\left(\dfrac{\partial x_j}{\partial x_k}\right)}_{\delta_{jk}} \right]$

$- \sum_i b_i \underbrace{\left(\dfrac{\partial x_i}{\partial x_k}\right)}_{\rightarrow \delta_{ik}}$

$\left[ \delta_{ik} = \begin{cases} 1 & \text{if } i = k \\ 0 & \text{else} \end{cases} \right]$

$= \dfrac{1}{2}\left( \sum_j A_{kj} x_j + \sum_i A_{ik} x_i \right)$

$\qquad - b_k$

$= \dfrac{1}{2}\left( \sum_i A_{ki} x_i + \sum_i A_{ik} x_i \right)$

$\qquad - b_k$

$= \sum_i A_{ki} x_i - b_k = 0$

$\Longleftrightarrow \quad A\underset{\sim}{x} = \underset{\sim}{b}$ $\left.\rule{0pt}{60pt}\right)$

② Solve min $Q$ using by an "iterative" optimization algorithm
(such as steepest descent method)

③ To speed up the convergence, may tranform $\underline{x}$ as

$$\underline{x} = \underline{\underline{T}} \, \underline{y}$$

with $\underline{\underline{T}} = [D_{ii}] \leftarrow$ Diagonal matrix

$$D_{ii} = \frac{1}{\sqrt{|A_{ii}|}}$$

## Numerical Example of the Steepest descent Method

Consider $f = (x_1-2)^4 + (x_1-2x_2)^2$, $x_0 = (0,3)^T$

$$f(x_0) = 52 \ , \ \nabla f(x) = \left[ 4(x_1 - 2)^3 + 2(x_1 - 2x_2), \ -4(x_1 - 2x_2) \right]^T$$

- Thus, the search direction is

$$d_0 = -\nabla f(x_0) = \left[ 44, -24 \right]^T$$

$$\underset{normalize}{\Rightarrow} \quad d_0 = \left[ 0.8779, -0.4789 \right]^T$$

- Line search

Minimize $f(\alpha) = f(x_0 + \alpha \, d_0)$ with $\alpha > 0$

$$\Rightarrow \alpha = 3.0841$$

$\therefore$ The new point is

$$x_1 = x_0 + \alpha_0 d_0 = [2.707, \ 1.523]^T$$