$$F(x + \Delta x) = F(x) + \Delta x \cdot \frac{dF}{dx} + \frac{1}{2}(\Delta x)^2 \frac{d^2F}{dx^2} + \cdots$$

Single Variable version.

$$H_{ij} = \frac{\partial^2 F}{\partial x_i \partial x_j}$$

If $X = (x_1, \cdots, x_m)$.

Hessian

$$F(x + \Delta x) = F(x) + (\Delta x)^T \nabla F(x) + \frac{1}{2}(\Delta x)^T H (\Delta x)$$

$$\nabla F(x) = \begin{bmatrix} \frac{\partial F}{\partial x_1} \\ \vdots \\ \frac{\partial F}{\partial x_m} \end{bmatrix}$$

$$[\cdots] \begin{bmatrix} \\ \\ \end{bmatrix}$$

$$H = \begin{bmatrix} \frac{\partial^2 F}{\partial x_1 \partial x_1} & \frac{\partial^2 F}{\partial x_1 \partial x_2} & \cdots \\ \frac{\partial^2 F}{\partial x_2 \partial x_1} & \ddots & \\ \vdots & & \frac{\partial^2 F}{\partial x_m \partial x_m} \end{bmatrix}$$
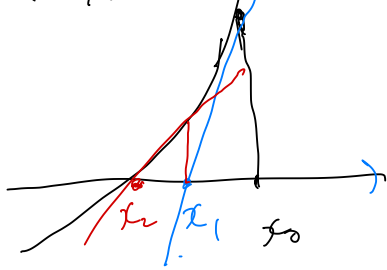
$$F(x,y) = 3x^2 + 2xy + 3y^2 \quad \Rightarrow \quad [x \quad y] \begin{bmatrix} 3 & 1 \\ 1 & 5 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}$$

$$f_x = 6x + 2y$$

$$f_{xx} = 6$$

$$H = \begin{bmatrix} f_{xx} & f_{xy} \\ f_{yx} & f_{yy} \end{bmatrix}$$

# Newton's Method (solve $f(x)=0$)



$$f(x_k + \Delta x) = f(x_k) + f'(x_k) \cdot \Delta x$$

$$\underbrace{\quad\quad\quad\quad}_{>0}$$

$$(\Delta x = x_{k+1} - x_k)$$

$$0 = f(x_k) + f'(x_k)(x_{k+1} - x_k)$$

$$\Rightarrow \boxed{x_{k+1} = x_k - \left(f'(x_k)\right)^{-1} \cdot f(x_k)}$$

$$\left( x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)} \right)$$

$$f = (f_1, f_2, f_3)$$

$$\Rightarrow \boxed{x_{k+1} = x_k - J(x_k)^{-1} \cdot f(x_k)}$$

$$\left( \text{Jacobian} \quad J_{jk} = \frac{\partial f_j}{\partial x_k} \right)$$

ex) $f(x) = x^2 - 9$.

$$x_{k+1} = x_k - J(x_k)^{-1} \cdot f(x_k).$$

$$J(x_k) = 2x_k \quad , \quad f(x_k) = x_k^2 - 9.$$

$$x_{k+1} = x_k - \frac{1}{2x_k}(x_k^2 - 9).$$

Convergence rate?

$$(x_{k+1} - 3) = x_k - \frac{1}{2x_k}(x_k^2 - 9) - 3.$$

$$= \frac{1}{2x_k}\{2x_k^2 - x_k^2 + 9 - 6x_k\}.$$

$$= \frac{1}{2x_k}(x_k - 3)^2.$$

$$\Rightarrow (x_{k+1} - 3) = \frac{1}{2x_k}(x_k - 3)^2.$$

→ very powerful method.

# Minimize F(x).    ($\approx$ solving $\nabla f = 0$).

(I) Steepest Descent.

$$X_{k+1} = X_k - S_F \cdot \nabla F$$

(II) Newton's Method.

$$\left( X_{k+1} - X_k - J^{-1}(x_k) \cdot f(x_k) \right) \qquad \text{now } f \to \nabla f)$$

$$J_{jk} = \frac{\partial f_j}{\partial x_k}. \longrightarrow \quad \nabla f_j = \frac{\partial f}{\partial x_j}. \quad \to \quad J_{jk} = \frac{\partial^2 f}{\partial x_k \partial x_j}$$

$$x_{k+1} = x_k - H^{-1}(\nabla F).$$

Convergence rate.

$\leftarrow$ linear

$\leftarrow$ quadratic.

cost?

$\leftarrow$ cheap.

$\leftarrow$ expensive (H).

# Stochastic Gradient descent. (SGD).

data 1 $\longrightarrow$    $w.$

data 2 $\longrightarrow$    [ NN ]

$\vdots$ $\longrightarrow$

data n $\longrightarrow$

$\rightarrow l_1 \rightarrow \nabla l_1$
$\rightarrow l_2 \rightarrow \nabla l_2$
$\vdots$
$\searrow l_n \rightarrow \nabla l_n$

$$\cancel{\nabla L = \frac{1}{n} \sum_{n=1}^{n} \nabla l_i}$$

$$\nabla L = \nabla l_k.$$

How pick data $k$? (for $k = 0, 1, \cdots, n$).

Option 1: Randomly pick an index $i$ with replacement.

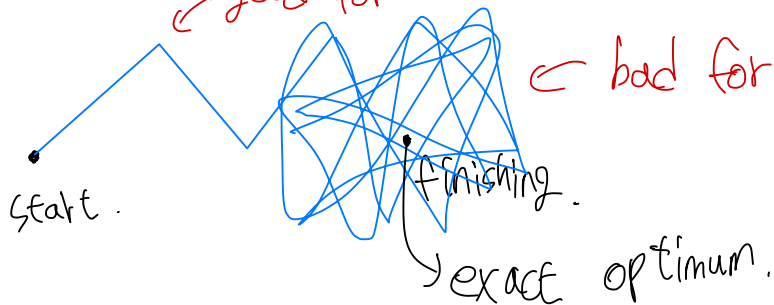Option 2: Pick index $i$ without replacement.

$\Rightarrow$ option 2!   ( shuffle 후   차례대로).

# SGD

Property

← good for beginning.



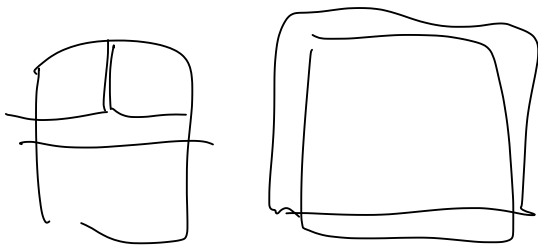← bad for finishing.

start.

finishing.

→ exact optimum.

If I don't care about getting to the best optimum?

⇒ "SGD is great"

(exact optimal → over fitting)

(Mini-batch).

→ GPU 연산에 적합.

Back-Propagation.

° To compute $\nabla f = \left( \frac{\partial F}{\partial x_1}, \cdots, \frac{\partial F}{\partial x_m} \right)$.

$\Rightarrow$ AD (automatic differentiation) "Reverse Mode".

순서가 중요?

$$A \underset{m \times n}{B} \underset{n \times p}{C} \underset{p \times q}{\rightarrow} \underbrace{\boxed{(AB)}}_{m \times p} \underset{p \times q}{C} = A \underset{m \times n}{(BC)} \underset{n \times q}{}$$

$$mnp + mpq \mid mnq + npq.$$

If $C = \begin{bmatrix} \ \end{bmatrix}$? $(q=1)$.

$$mnp + mp \ggg mn + np.$$

$$\begin{bmatrix} \oslash\oslash & \oslash \end{bmatrix} \begin{bmatrix} \ominus \\ \ominus \\ \ominus \end{bmatrix} = \begin{bmatrix} \bigcirc & \cdot \\ \cdot & \end{bmatrix}$$

$\underline{2 \times 3}$   $\underline{3 \times 2}$   $\underline{2 \times 2}$.

# In forward - mode



$$\frac{\partial e}{\partial x_{4}} \qquad \frac{\partial e}{\partial x_{1}}$$

$$\frac{\partial e}{\partial x_{2}}.$$

$$\frac{\partial e}{\partial x_{3}}$$

# In backward - mode.



$$\frac{\partial W_{3}}{\partial W_{6}} \qquad \frac{\partial W_{1}}{\partial W_{3}} \qquad \frac{\partial e}{\partial W_{1}}.$$

$$W_{6} \qquad W_{3}$$

$$W_{1}$$

$$\frac{\partial W_{4}}{\partial W_{7}} \qquad W_{4} \frac{\partial W_{1}}{\partial W_{4}}.$$

$$W_{4} \qquad W_{5} \qquad W_{2} \frac{\partial e}{\partial W_{2}}.$$

$$\frac{\partial W_{5}}{\partial W_{7}} \qquad \frac{\partial W_{2}}{\partial W_{5}}$$

$$\frac{\partial e}{\partial W_{6}} = \frac{\partial e}{\partial W_{1}} \cdot \frac{\partial W_{1}}{\partial W_{3}} \cdot \frac{\partial W_{3}}{\partial W_{6}}.$$

$$\frac{\partial e}{\partial W_{7}} = \frac{\partial e}{\partial W_{1}} \cdot \frac{\partial W_{1}}{\partial W_{4}} \cdot \frac{\partial W_{4}}{\partial W_{7}} + \frac{\partial e}{\partial W_{2}} \cdot \frac{\partial W_{2}}{\partial W_{5}} \cdot \frac{\partial W_{5}}{\partial W_{7}}.$$

"ten million faster"    1000만배.