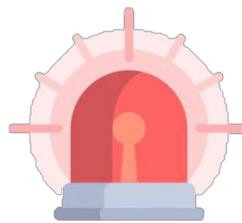


# Regression with Gradient Descent Method



Complex: Differentiation, Derivative, and Optimization

One of difficult parts in this class, but  
it will be a foundation of Deep Neural Networks!!

# Regression

- ▶ Solve the following problem

*Given Data* =  $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$

Model you choose  $f(\mathbf{x}; \mathbf{w})$

Find  $\mathbf{w}$  to minimizes  $E$  ( $E_{\text{rror}}$ )

$$E(\mathbf{w}) = \sum_{(\mathbf{x}, y) \in \text{Data}} (y - \underline{f(\mathbf{x}; \mathbf{w})})^2$$

Regression is an Optimization problem !!

# Linear Additive Model

- ▶ E is a quadratic function of w

Find  $w_1, w_2, \dots, w_m$  to minimize

$$E(w_1, w_2, \dots, w_m) = \sum_{(\mathbf{x}, y) \in Data} (y - f(\mathbf{x}; w_1, w_2, \dots, w_m))^2$$

f is Linear fun of w

- ▶ It is a quadratic function optimization. <- Easy

$$\left. \begin{array}{l} \frac{\partial E}{\partial w_1} = 0 \\ \frac{\partial E}{\partial w_2} = 0 \\ \dots \\ \frac{\partial E}{\partial w_m} = 0 \end{array} \right\}$$

Solve the linear equations. Then we have the exact solution !!

# But..

$$E = -2 \left( y - \underbrace{w_0 e^{wx} - \sin wx}_{f(x)} \right) (we^{wx} + w \cos wx) = 0 \quad \therefore \text{부정해석 유풍..}$$

$\hookrightarrow \frac{dE}{dw} = 0$

- ▶ What if  $f$  is not quadratic nor a polynomial of  $w$

$$\frac{\partial E}{\partial w_1} = 0$$

$$\frac{\partial E}{\partial w_2} = 0$$

...

$$\frac{\partial E}{\partial w_m} = 0$$

Can we solve the equations?

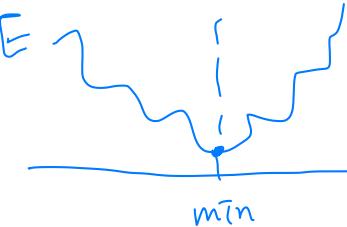
일반적으로, 찾기 어렵다 (계산이 복잡해지기)

값을 알면 어렵지만 분명히 존재함

- ▶ Is there a solution?

Yes!

$$E(w_1, w_2, \dots, w_m) = \sum_{(\mathbf{x}, y) \in Data} (y - f(\mathbf{x}; w_1, w_2, \dots, w_m))^2$$



- ▶ Can we find the solution?
- ▶ What shall we do? :(

We canNOT say "Yes"

# General Approach for Optimization

- Because the original problem ~~can~~ may not be solved

Find  $w$  that (globally) minimizes

$$E(w) = \sum_{(x,y) \in Data} (y - f(x; w))^2$$

non-linear

no solution / (fail)

$\frac{dE}{dw} = 0$  local min

- We change the problem as follows:

Find  $w$  that locally minimizes

$$E(w) = \sum_{(x,y) \in Data} (y - f(x; w))^2$$

Not Exact Solution  
by GDM

Local Minima!

# General Approach for Optimization

## ▶ Local Minimization instead of Global Minimization

Find  $w$  that locally minimizes

$$E(w) = \sum_{(x,y) \in Data} (y - f(x; w))^2$$

repeat

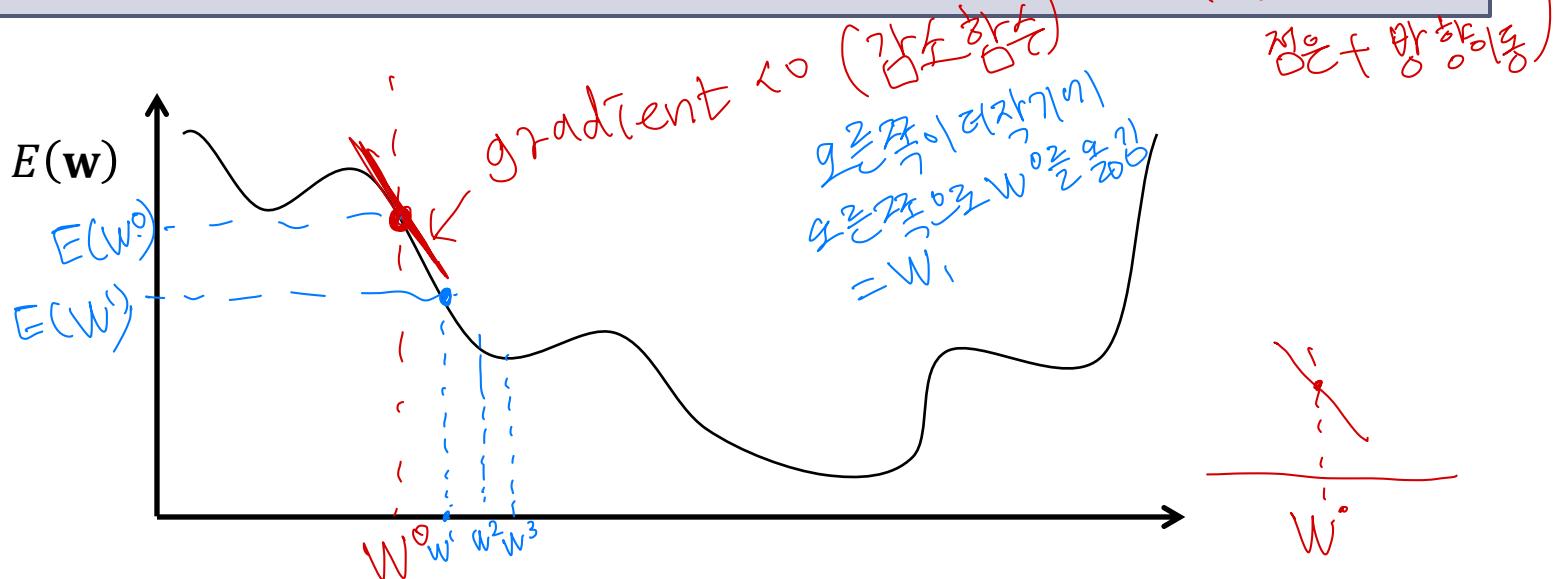
① Random  $w^0$

②  $\frac{\partial E}{\partial w}$  at  $w^0$

③ move opposite direction

(기울기가 -면)

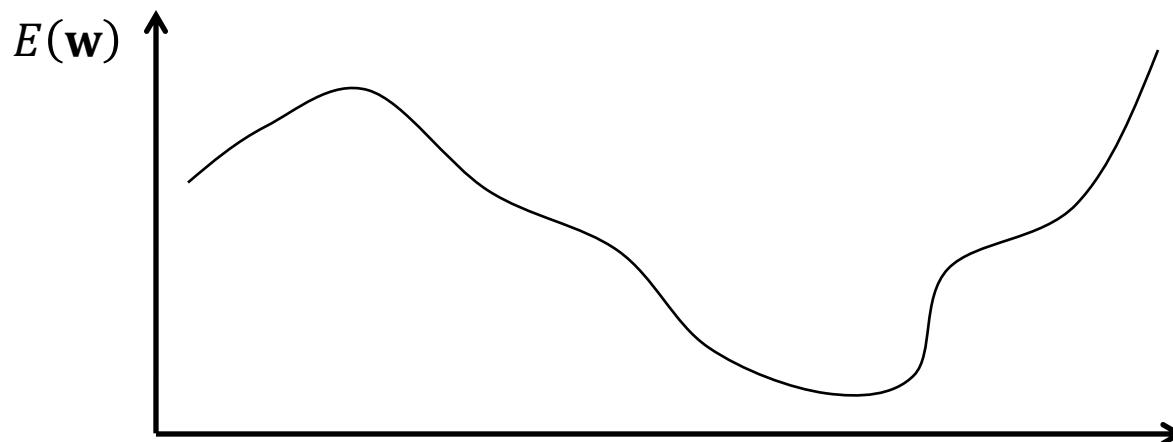
점은 + 방향이동



# General Approach for Optimization

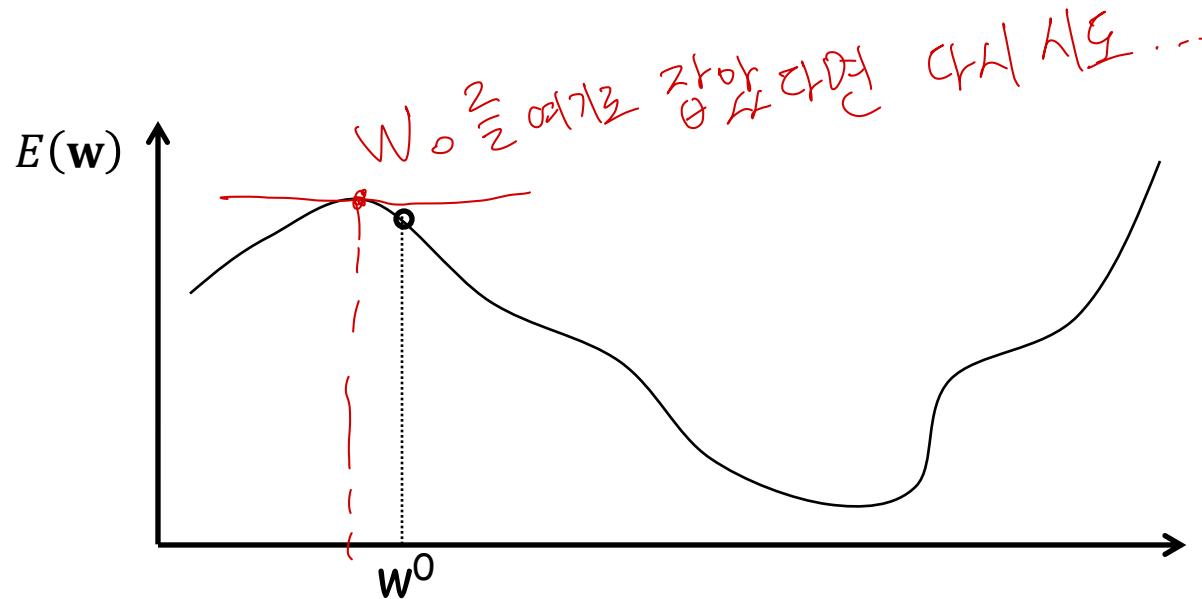
- ▶ OK, but.. Is it easy to find a local minimum? Yes!!

But, how?



# Gradient Descent Method

- ▶ 1. Randomly choose an initial point



# Gradient Descent Method

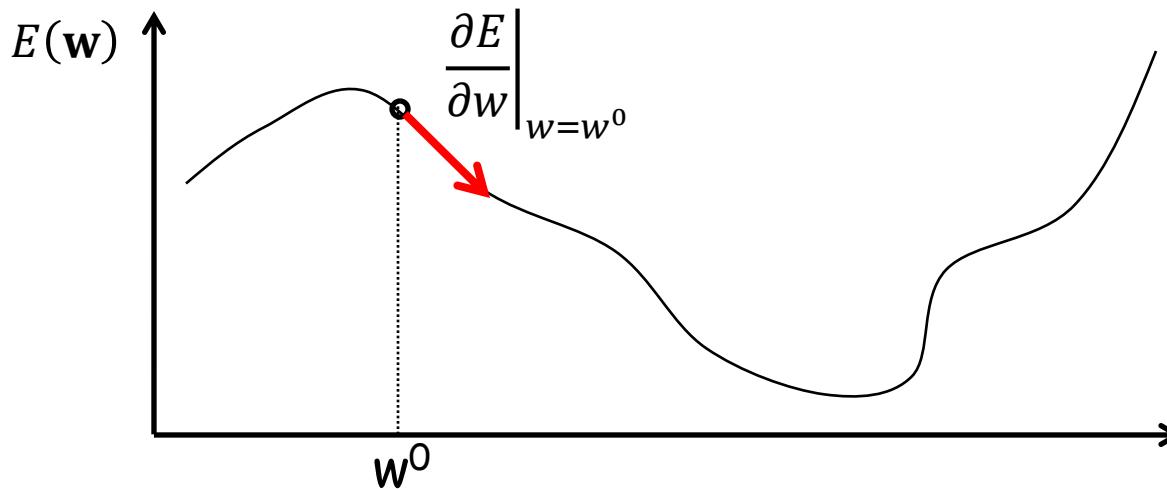
## ▶ 2. Evaluate the gradient

▶ E(w) is differentiable -> We can obtain the gradient

assumption

$$\frac{\partial E}{\partial w}$$

이alon가능한

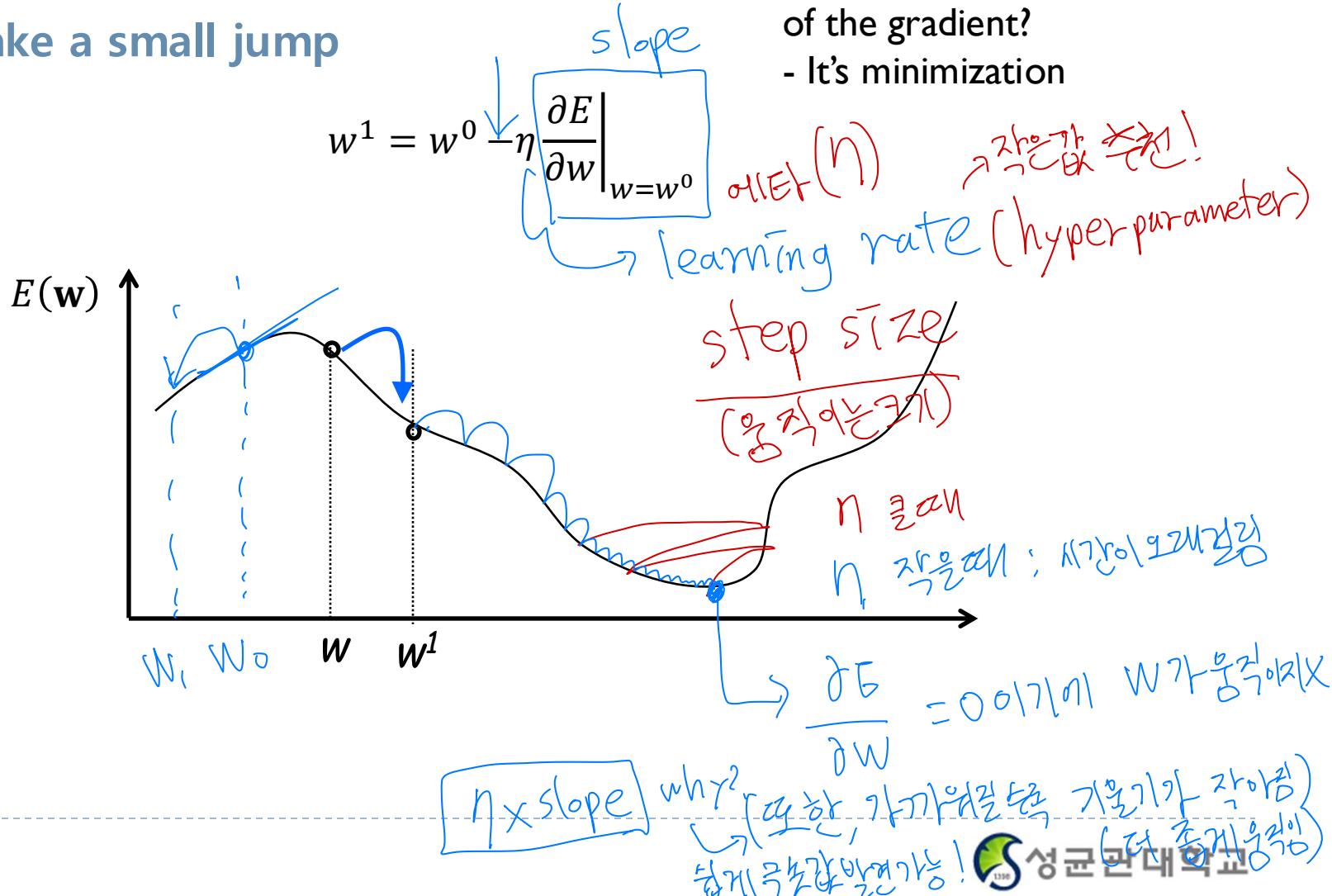


# Gradient Descent Method

- ▶ 3. Move downward
    - ▶ Make a small jump

→ 느리기이  
**d** 각 속도  $v$  여러 개를 찾기는 힘들다  
∴ local minimum 1개만 찾는다

Why to jump reverse direction  
of the gradient?  
- It's minimization

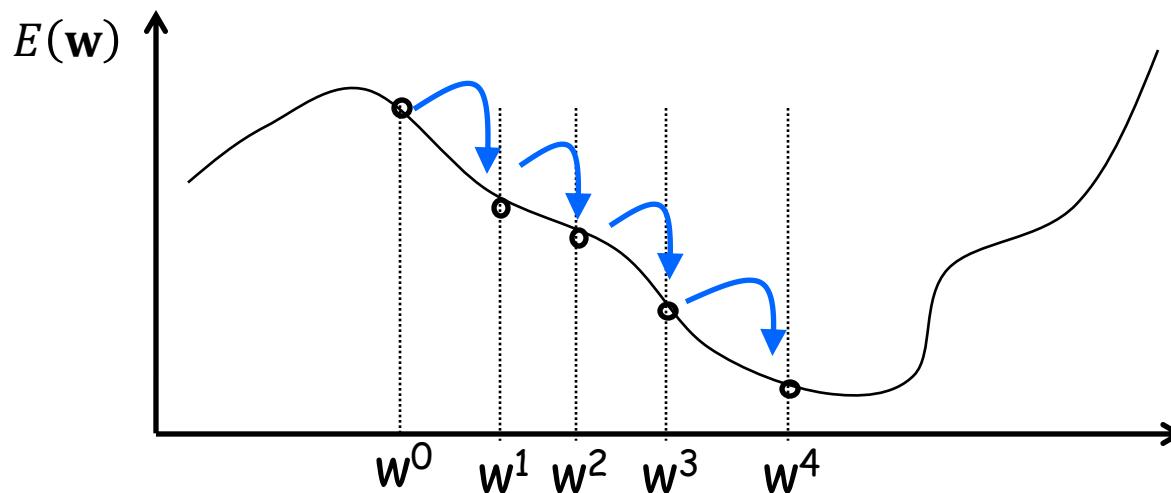


# Gradient Descent Method

- ▶ 4. Repeat the steps until the gradient is zero

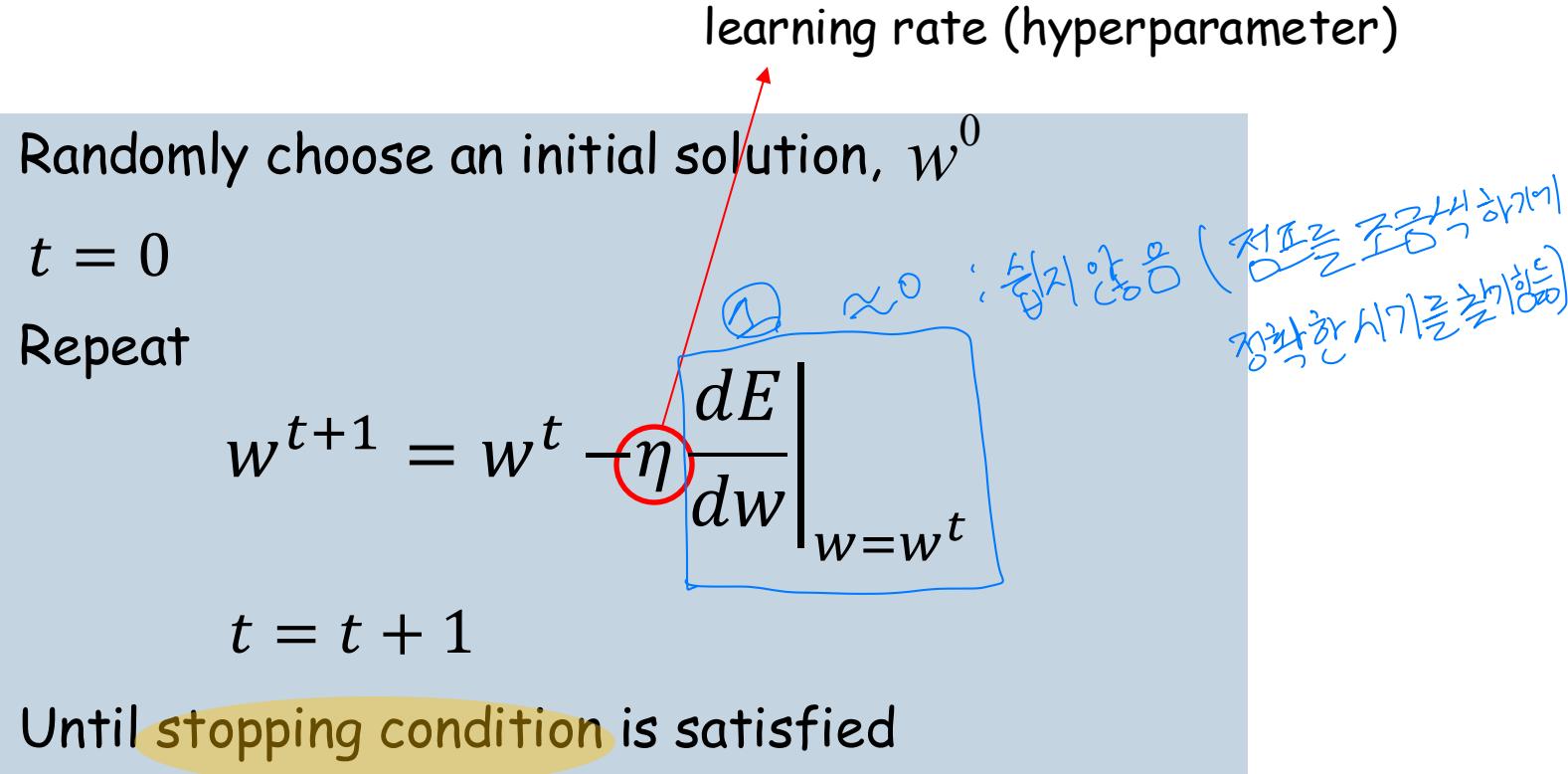
adam optimizer  
L.r이 고정된 것.  
"adaptive" 시기  
설정하는 것의 종료

$$w^{t+1} = w^t - \eta \left. \frac{\partial E}{\partial w} \right|_{w=w^t}$$



# Gradient Descent Method

## ▶ Steps



# Gradient Descent Method

시간복잡도 높다

$$\frac{1}{2}w^2 + w + 1$$

## Steps

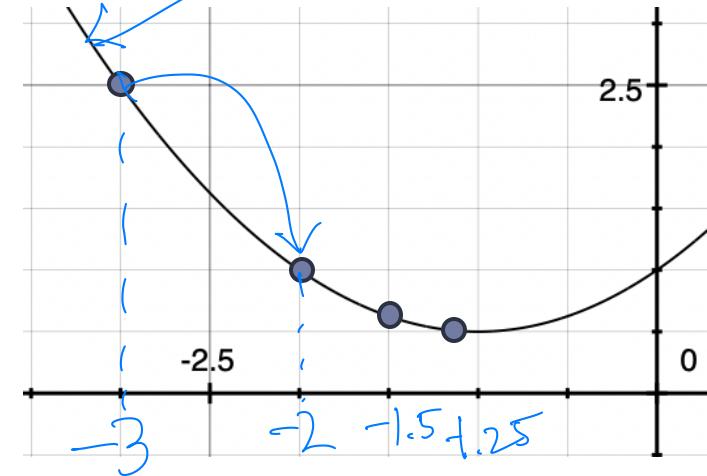
Randomly choose an initial solution,  $w^0$

Repeat

$$w^{t+1} = w^t - \eta \frac{dE}{dw} \Big|_{w=w^t}$$

$t = t + 1$

Until stopping condition is satisfied



$$E = \frac{1}{2}w^2 + w + 1$$

$$w^0 = -3$$

$$\frac{dE}{dw} = w + 1$$

$$\eta = 0.5$$

$$w^1 = w^0 - \eta \frac{dE}{dw} \Big|_{w=w^0} = -3 - 0.5 \times (-2) = -2$$

$$w^2 = w^1 - \eta \frac{dE}{dw} \Big|_{w=w^1} = -2 - 0.5 \times (-1) = -1.5$$

$$w^3 = w^2 - \eta \frac{dE}{dw} \Big|_{w=w^2} = -1.5 - 0.5 \times (-0.5) = -1.25$$

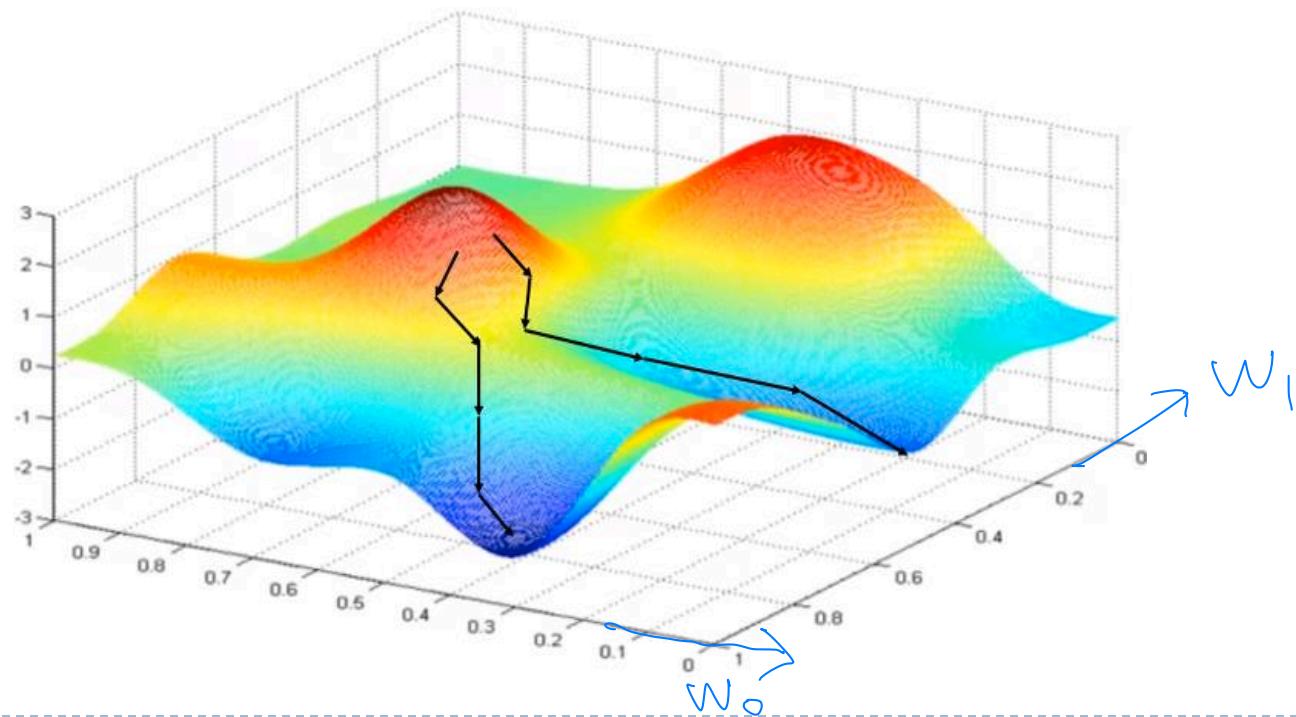
next point



# Gradient Descent Method

## ▶ Multivariate Case

$$E(\mathbf{w}) = \sum_{(\mathbf{x}, y) \in Data} (y - f(\mathbf{x}; \mathbf{w}))^2$$



# Gradient Descent Method

## Multivariate Case

- Separately apply the method to each variable

$w_0^0$   $\nwarrow$  iteration  
 $w_0^0$   $\nwarrow$  variable  
initial point

Randomly choose an initial solution,  $w_0^0 w_1^0$

$t = 0$

2-dimension

Repeat

$$w_0^{t+1} = w_0^t - \eta \frac{\partial E}{\partial w_0} \quad \left|_{w_0=w_0^t, w_1=w_1^t} \right.$$
$$w_1^{t+1} = w_1^t - \eta \frac{\partial E}{\partial w_1} \quad \left|_{w_0=w_0^t, w_1=w_1^t} \right.$$

기울기  
이동  
 $\Rightarrow$  그대로 광학적  
현재 포지션

$t = t + 1$

Until stopping condition is satisfied

# Gradient Descent Method

## ▶ Steps

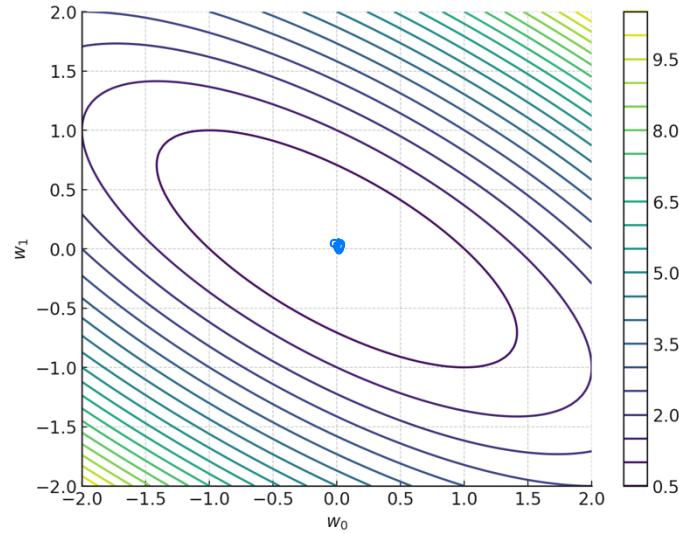
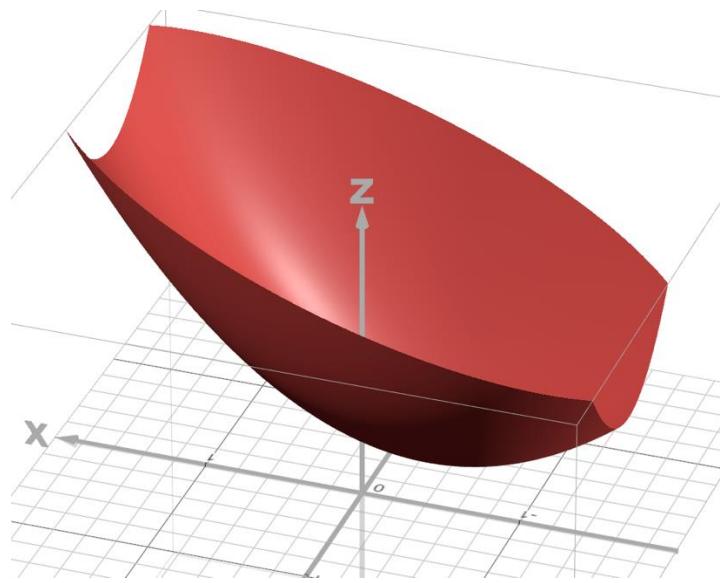
$$E = \frac{1}{2}w_0^2 + w_1^2 + w_0w_1 + \frac{1}{2}$$

$w_0=0 \quad w_1=0$

optimization

$$\frac{\partial E}{\partial w_0} = w_0 + w_1$$

$$\frac{\partial E}{\partial w_1} = 2w_1 + w_0$$



# Gradient Descent Method

## ▶ Steps

Randomly choose an initial solution,  $w_0^0 w_1^0$   
 $t = 0$   
 Repeat  
 $w_0^{t+1} = w_0^t - \eta \frac{\partial E}{\partial w_0} \Big|_{w_0=w_0^t, w_1=w_1^t}$   
 $w_1^{t+1} = w_1^t - \eta \frac{\partial E}{\partial w_1} \Big|_{w_0=w_0^t, w_1=w_1^t}$   
 $t = t + 1$   
 Until stopping condition is satisfied

$$E = \frac{1}{2} w_0^2 + w_1^2 + w_0 w_1 + \frac{1}{2}$$

$$\frac{\partial E}{\partial w_0} = w_0 + w_1 = -3$$

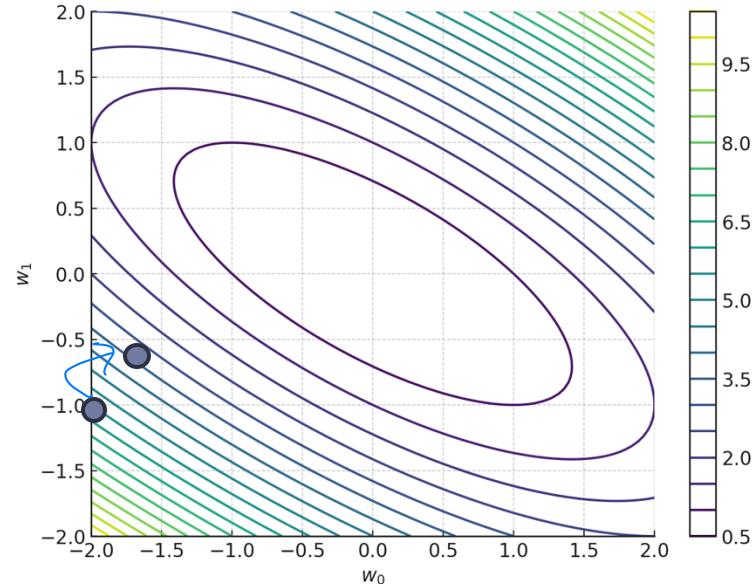
$$w_0^0 = -2 \quad w_1^0 = -1$$

$$\frac{\partial E}{\partial w_1} = 2w_1 + w_0 = -4$$

$$\eta = 0.1$$

$$w_0^1 = w_0^0 - \eta \frac{\partial E}{\partial w_0} \Big|_{w_0=w_0^0, w_1=w_1^0} = -2 - 0.1 \times (-3) = -1.7$$

$$w_1^1 = w_1^0 - \eta \frac{\partial E}{\partial w_1} \Big|_{w_0=w_0^0, w_1=w_1^0} = -1 - 0.1 \times (-4) = -0.6$$



# Gradient Descent Method

## ▶ Steps

Randomly choose an initial solution,  $w_0^0 w_1^0$

$t = 0$

Repeat

$$w_0^{t+1} = w_0^t - \eta \frac{\partial E}{\partial w_0} \Big|_{w_0=w_0^t, w_1=w_1^t}$$

$$w_1^{t+1} = w_1^t - \eta \frac{\partial E}{\partial w_1} \Big|_{w_0=w_0^t, w_1=w_1^t}$$

$t = t + 1$

Until stopping condition is satisfied

$$E = \frac{1}{2} w_0^2 + w_1^2 + w_0 w_1 + \frac{1}{2}$$

$$\frac{\partial E}{\partial w_0} = w_0 + w_1$$

$$w_0^1 = -1.7 \quad w_1^1 = -0.6$$

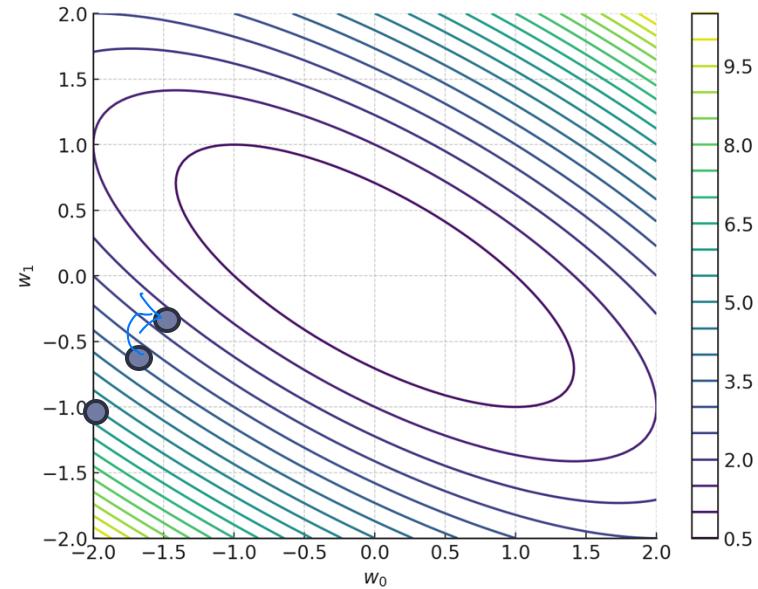
$$w_0^2 = -1.47 \quad w_1^2 = -0.31$$

$$\frac{\partial E}{\partial w_1} = 2w_1 + w_0$$

$$\eta = 0.1$$

$$w_0^2 = w_0^1 - \eta \frac{\partial E}{\partial w_0} \Big|_{w_0=w_0^1, w_1=w_1^1} = -1.7 - 0.1 \times (-2.3) = -1.47$$

$$w_1^2 = w_1^1 - \eta \frac{\partial E}{\partial w_1} \Big|_{w_0=w_0^1, w_1=w_1^1} = -0.6 - 0.1 \times (-2.9) = -0.31$$



# Gradient Descent Method

## ▶ Steps

Randomly choose an initial solution,  $w_0^0 w_1^0$

$t = 0$

Repeat

$$w_0^{t+1} = w_0^t - \eta \frac{\partial E}{\partial w_0} \Big|_{w_0=w_0^t, w_1=w_1^t}$$

$$w_1^{t+1} = w_1^t - \eta \frac{\partial E}{\partial w_1} \Big|_{w_0=w_0^t, w_1=w_1^t}$$

$t = t + 1$

Until stopping condition is satisfied

$$E = \frac{1}{2} w_0^2 + w_1^2 + w_0 w_1 + \frac{1}{2}$$

$$\frac{\partial E}{\partial w_0} = w_0 + w_1$$

$$w_0^2 = -1.47 \quad w_1^2 = -0.31$$

$$w_0^2 = -1.292 \quad w_1^2 = -0.101$$

$$\frac{\partial E}{\partial w_1} = 2w_1 + w_0$$

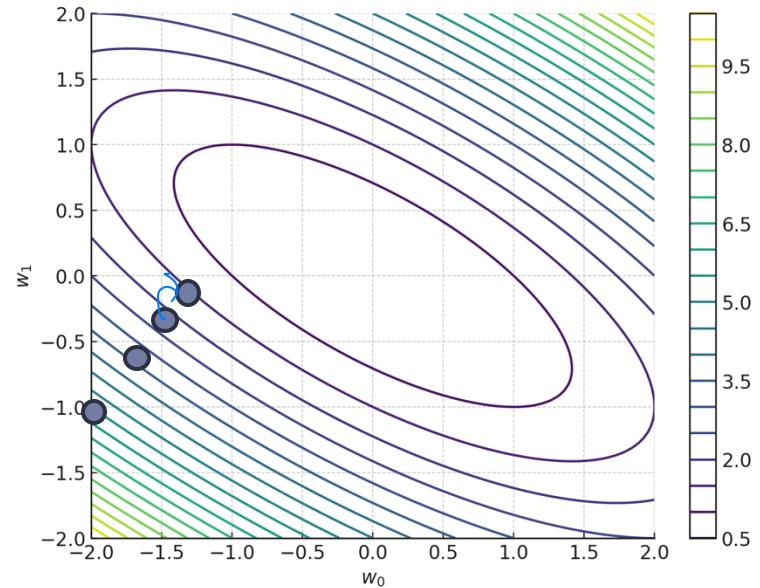
$$\eta = 0.1$$

$$w_0^3 = w_0^2 - \eta \frac{\partial E}{\partial w_0} \Big|_{w_0=w_0^2, w_1=w_1^2}$$

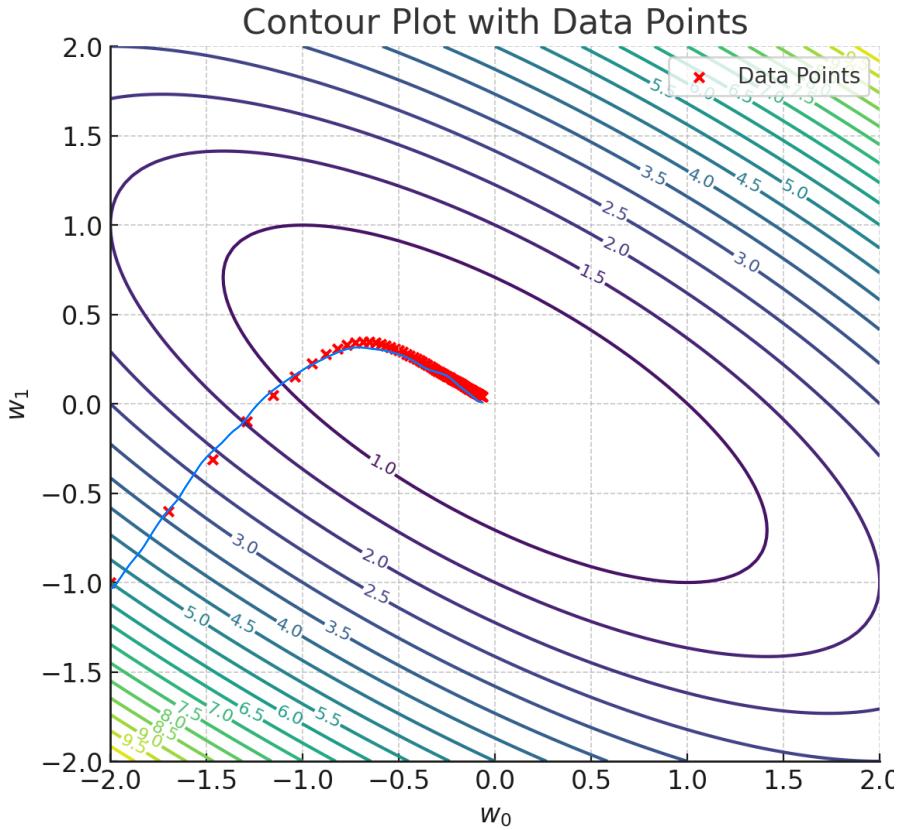
$$w_1^3 = w_1^2 - \eta \frac{\partial E}{\partial w_1} \Big|_{w_0=w_0^2, w_1=w_1^2}$$

$$= -1.47 - 0.1 \times (-1.78) = -1.292$$

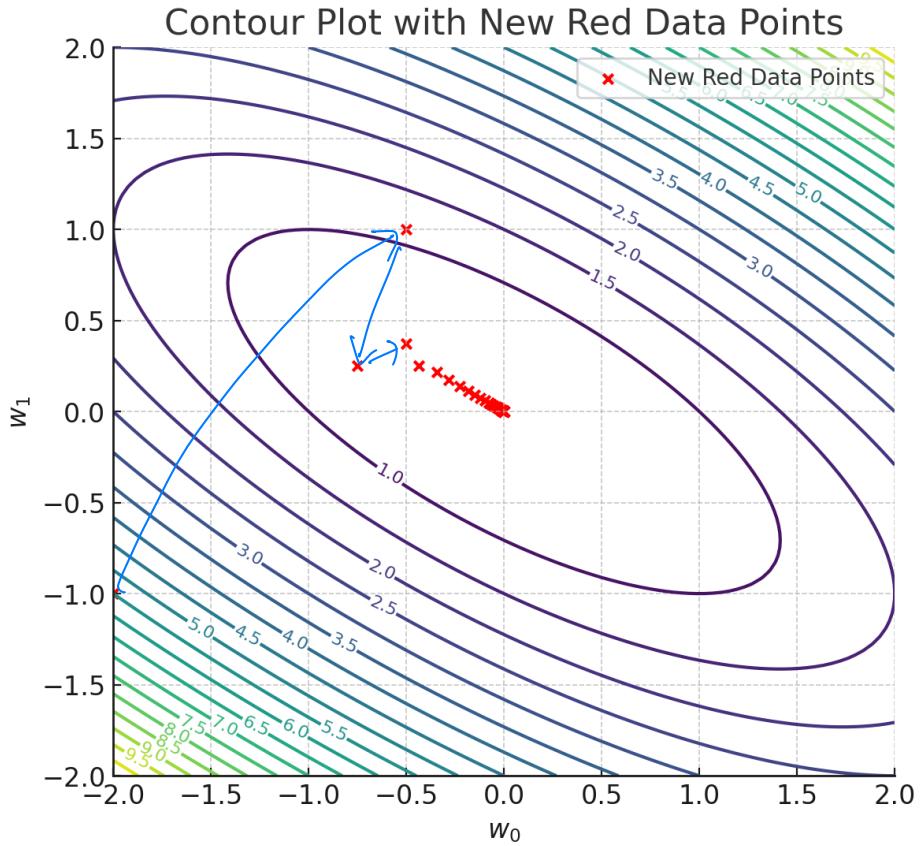
$$= -0.31 - 0.1 \times (-2.09) = -0.101$$



# Gradient Descent Method



$$\eta = 0.1$$



$$\eta = 0.5$$

Large  $\eta$   
Slow Convergence  
Large Step Size  
Large Step Size

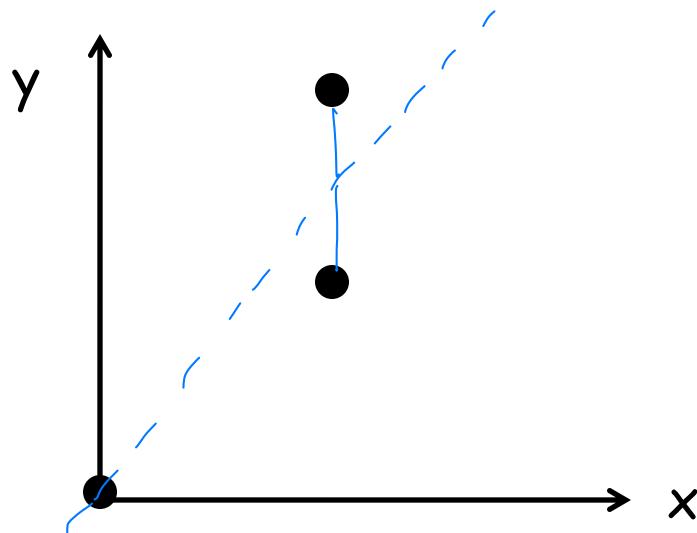
성균관대학교

# **Linear Regression with Gradient Descent**

# Linear Regression

- ▶ Find the best-fit line

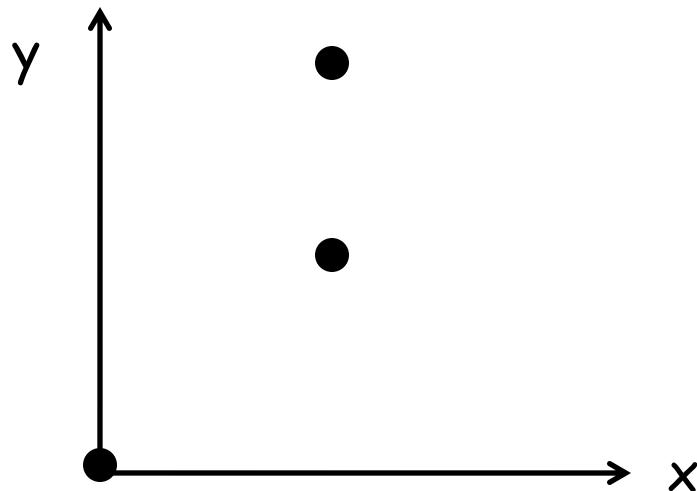
$(0.0, 0.0) (1.0, 1.0) (1.0, 2.0)$



# Linear Regression

- ▶ Step 1: Choose a model

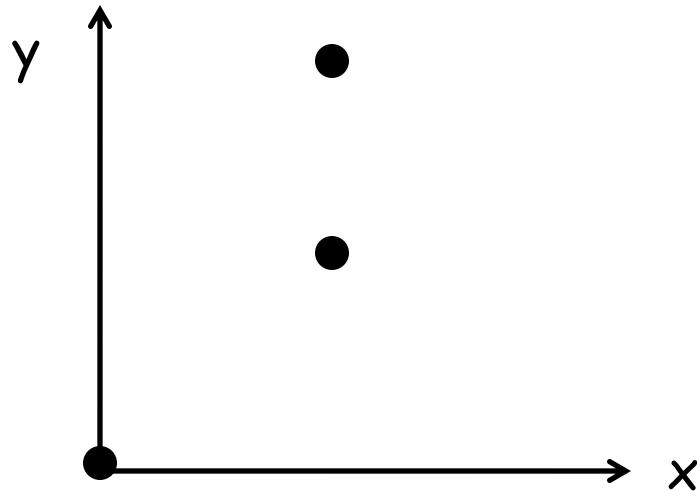
$$f(x; w_0, w_1) = w_1 x + w_0$$



# Steps of Machine Learning

- ▶ Step 2: Find  $w$  to minimize  $E$  by GDM

$$f(x; w_0, w_1) = w_1 x + w_0$$



# Steps of Machine Learning

## ▶ Step 2-1: Obtain a closed form of E

### ▶ Error Function

$$E(w_0, w_1) = \sum_{(x,y) \in Data} (y - f(x; w_1, w_0))^2$$

### ▶ Obtain a closed form of E by applying the given data

$$(0.0, 0.0) (1.0, 1.0) (1.0, 2.0) \quad f(x; w_0, w_1) = w_1 x + w_0$$

$$E(w_0, w_1) = (0.0 - f(0.0; w_0, w_1))^2 + (1.0 - f(1.0; w_0, w_1))^2 + (2.0 - f(1.0; w_0, w_1))^2$$

$$E(w_0, w_1) = (0.0 - w_0)^2 + (1.0 - (w_1 + w_0))^2 + (2.0 - (w_1 + w_0))^2$$

$$E(w_0, w_1) = 2w_1^2 + 3w_0^2 - 6w_1 - 6w_0 + 4w_1w_0 + 5$$

# Steps of Machine Learning

## ▶ Step 2-2: Obtain partial derivatives of E

error fun



$$E(w_0, w_1) = 2w_1^2 + 3w_0^2 - 6w_1 - 6w_0 + 4w_1w_0 + 5$$

E는 학습하고자 하는  
w<sub>0</sub>, w<sub>1</sub>에 대한

= optimization  
problem

≠ machine  
learning  
problem

$$\frac{\partial E}{\partial w_0} = 4w_1 + 6w_0 - 6$$

$$\frac{\partial E}{\partial w_1} = 4w_1 + 4w_0 - 6$$

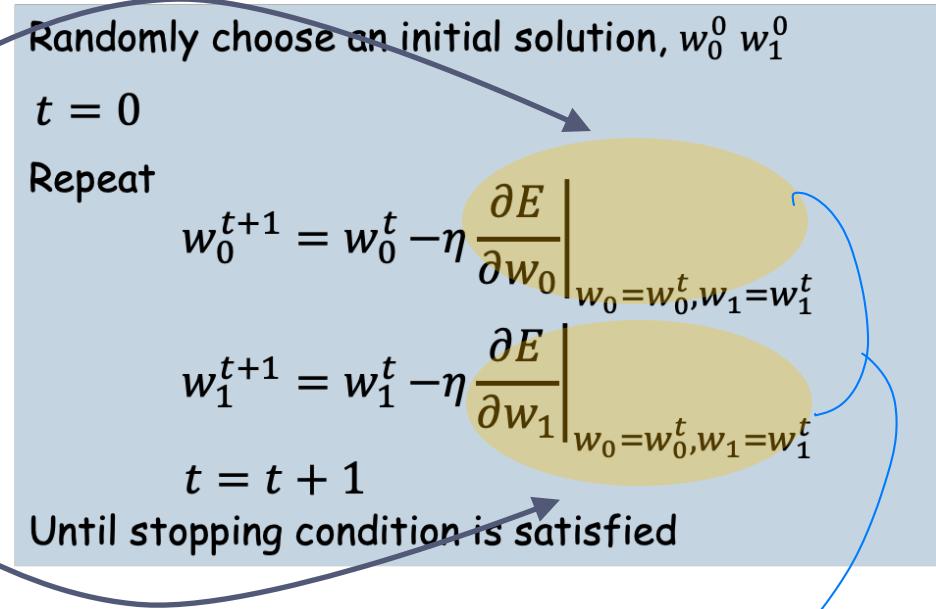
# Steps of Machine Learning

## ▶ Step 2-3: Plug the derivatives into the algorithm

$$\frac{\partial E}{\partial w_0} = 4w_1 + 6w_0 - 6$$

$$\frac{\partial E}{\partial w_1} = 4w_1 + 4w_0 - 6$$

### ▶ Final algorithm



Randomly choose an initial solution,  $w_0^0 w_1^0$

$t = 0$

Repeat

$$w_0^{t+1} = w_0^t - \eta(4w_1^t + 6w_0^t - 6)$$

$$w_1^{t+1} = w_1^t - \eta(4w_1^t + 4w_0^t - 6)$$

$t = t + 1$

Until stopping condition is satisfied

# Steps of Machine Learning

## ▶ Step 2-4: Run the algorithm

- ▶ Choose a small value for  $\eta$  ( $\eta = 0.1$ )
- ▶ Randomly choose  $w_0^0, w_1^0$ . ( $w_0^0 = 1, w_1^0 = 1$ )

$$\begin{aligned}w_0^0 &= 1 \\w_1^0 &= 1\end{aligned}$$

$$\begin{aligned}w_0^1 &= 1 - 0.1(4 \times 1 + 6 \times 1 - 6) = 0.6 \\w_1^1 &= 1 - 0.1(4 \times 1 + 4 \times 1 - 6) = 0.8\end{aligned}$$

$$\begin{aligned}w_0^2 &= 0.6 - 0.1(4 \times 0.8 + 6 \times 0.6 - 6) = 0.54 \\w_1^2 &= 0.8 - 0.1(4 \times 0.8 + 4 \times 0.6 - 6) = 0.84\end{aligned}$$

$$\begin{aligned}w_0^3 &= 0.54 - 0.1(4 \times 0.84 + 6 \times 0.54 - 6) = 0.480 \\w_1^3 &= 0.84 - 0.1(4 \times 0.84 + 4 \times 0.54 - 6) = 0.888\end{aligned}$$

Randomly choose an initial solution,  $w_0^0 w_1^0$   
 $t = 0$   
Repeat  
 $w_0^{t+1} = w_0^t - \eta \frac{\partial E}{\partial w_0} \Big|_{w_0=w_0^t, w_1=w_1^t}$   
 $w_1^{t+1} = w_1^t - \eta \frac{\partial E}{\partial w_1} \Big|_{w_0=w_0^t, w_1=w_1^t}$   
 $t = t + 1$   
Until stopping condition is satisfied

# Steps of Machine Learning

## ▶ Step 2-4: Run the algorithm

$$w_0^3 = \textcolor{violet}{0.54} - 0.1(4 \times \textcolor{red}{0.84} + 6 \times \textcolor{violet}{0.54} - 6) = \textcolor{violet}{0.480}$$

$$w_1^3 = \textcolor{red}{0.84} - 0.1(4 \times \textcolor{red}{0.84} + 4 \times \textcolor{violet}{0.54} - 6) = \textcolor{red}{0.888}$$

$$w_0^4 = \textcolor{violet}{0.480} - 0.1(4 \times \textcolor{red}{0.888} + 6 \times \textcolor{violet}{0.480} - 6) = \textcolor{violet}{0.4368}$$

$$w_1^4 = \textcolor{red}{0.888} - 0.1(4 \times \textcolor{red}{0.888} + 4 \times \textcolor{violet}{0.480} - 6) = \textcolor{red}{0.9408}$$

...

$$w_0^{100} = \textcolor{violet}{0.00007713}$$

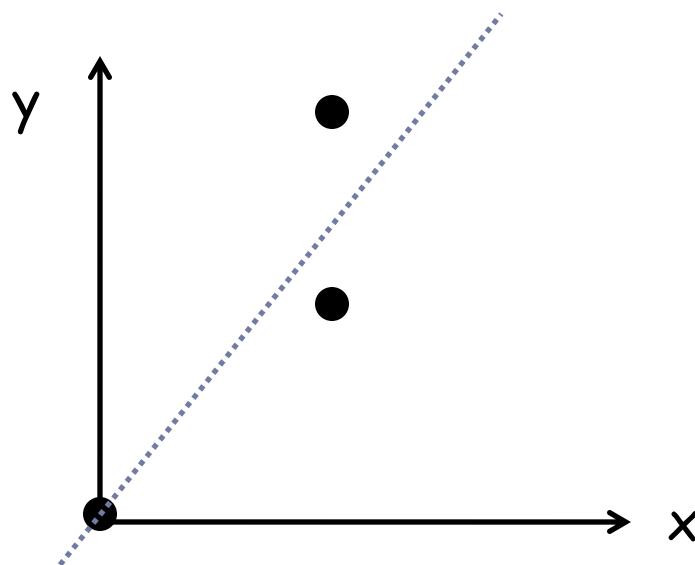
$$w_1^{100} = \textcolor{red}{1.49989171}$$

$w_0 \approx 0$   
 $w_1 \approx 1.5$

# Steps of Machine Learning

## ▶ Final Result

$$f(x) = 1.49989171x + 0.00007713$$



# But.. if I have millions of samples

- If there are a small number of samples, we can manually obtain derivatives

$$E(w_0, w_1) = \sum_{(\mathbf{x}, t) \in Data} (t - f(\mathbf{x}; w_1, w_0))^2$$

- For three samples, we manually obtain the closed form of E

$$E(w_0, w_1) = (\mathbf{0.0} - f(\mathbf{0.0}; w_0, w_1))^2 + (\mathbf{1.0} - f(\mathbf{1.0}; w_0, w_1))^2 + (\mathbf{2.0} - f(\mathbf{1.0}; w_0, w_1))^2$$

$$E(w_0, w_1) = (\mathbf{0.0} - w_0)^2 + (\mathbf{1.0} - (w_1 + w_0))^2 + (\mathbf{2.0} - (w_1 + w_0))^2$$

$$E(w_0, w_1) = 2w_1^2 + 3w_0^2 - 6w_1 - 6w_0 + 4w_1w_0 + 5$$

- We manually obtain derivatives

$$\frac{\partial E}{\partial w_0} = 4w_1 + 6w_0 - 6 \quad \frac{\partial E}{\partial w_1} = 4w_1 + 4w_0 - 6$$

- How can I MANUALLY obtain derivatives? if I have millions of samples.. :(

# But.. if I have millions of samples

definition

$$E(w_0, w_1) = \sum_{(\mathbf{x}, y) \in Data} (y - f(\mathbf{x}; w_1, w_0))^2$$

Closed form

of  $E$   
manually

$$\begin{aligned} E(w_0, w_1) = & 2w_1^2 + 3w_0^2 - 6w_1 \\ & - 6w_0 + 4w_1w_0 + 5 \end{aligned}$$

Differentiation

$$\frac{\partial E}{\partial w_0} = 4w_1 + 6w_0 - 6$$

$$\frac{\partial E}{\partial w_1} = 4w_1 + 4w_0 - 6$$

sample  $\mathbf{x}, y$  of  $E$

Differentiation without  
closed form of  $E$

closed form of  $E$  필요 없어  
directly differentiate

$$\frac{\partial E}{\partial w_0} = \sum_{(\mathbf{x}, y) \in Data} \frac{\partial}{\partial w_0} (y - f(\mathbf{x}; w_1, w_0))^2$$

$$\frac{\partial E}{\partial w_1} = \sum_{(\mathbf{x}, y) \in Data} \frac{\partial}{\partial w_1} (y - f(\mathbf{x}; w_1, w_0))^2$$

# But.. if I have millions of samples

- ▶ Another way to obtain derivatives at a point

$$E(w_0, w_1) = \sum_{(\mathbf{x}_i, t_i) \in Data} (t_i - f(\mathbf{x}_i; w_1, w_0))^2 = E_i \text{ (error of } \mathbf{x}_i)$$

*Sum of E*

$$= \sum_{(\mathbf{x}_i, t_i) \in Data} E_i(\mathbf{x}_i, t_i; w_0, w_1)$$

where  $E_i(\mathbf{x}_i, t_i; w_1, w_0) = (t_i - f(\mathbf{x}_i; w_1, w_0))^2$

▶ then

$$\frac{\partial}{\partial w_j} E(w_0, w_1) = \frac{\partial}{\partial w_j} \sum_{(\mathbf{x}_i, t_i) \in Data} E_i(\mathbf{x}_i, t_i; w_1, w_0)$$

$$= \sum_{(\mathbf{x}_i, t_i) \in Data} \frac{\partial}{\partial w_j} E_i(\mathbf{x}_i, t_i; w_1, w_0)$$

*j번째 derivative를 Σ*

*i → sample을 구현하는 단위  
j → sample을 구성하는 성분  
or  
w의 성분*

# But.. if I have millions of samples

- ▶ Another way to obtain derivatives

- ▶ Derivative of E

$$\frac{\partial E}{\partial w_j} = \sum_{(\mathbf{x}_i, t_i) \in Data} \frac{\partial E_i}{\partial w_j}$$

- ▶ Derivative of E at  $(w_0^t, w_1^t)$

$$\left. \frac{\partial E}{\partial w_j} \right|_{w_0=w_0^t, w_1=w_1^t} = \sum_{(\mathbf{x}_i, t_i) \in Data} \left. \frac{\partial E_i}{\partial w_j} \right|_{w_0=w_0^t, w_1=w_1^t}$$

전체 E = 각  $E_i$ 의  $\Sigma$

# But.. if I have millions of samples

## ▶ Another way to obtain derivatives

Randomly choose an initial solution,  $w_0^0 w_1^0$

$t = 0$

Repeat

$$w_0^{t+1} = w_0^t - \eta \frac{\partial E}{\partial w_0} \Big|_{w_0=w_0^t, w_1=w_1^t}$$

$$w_1^{t+1} = w_1^t - \eta \frac{\partial E}{\partial w_1} \Big|_{w_0=w_0^t, w_1=w_1^t}$$

$t = t + 1$

Until stopping condition is satisfied

$$\sum_{(\mathbf{x}_i, t_i) \in Data} \frac{\partial E_i}{\partial w_0} \Big|_{w_0=w_0^t, w_1=w_1^t}$$

$$\sum_{(\mathbf{x}_i, t_i) \in Data} \frac{\partial E_i}{\partial w_1} \Big|_{w_0=w_0^t, w_1=w_1^t}$$

# But.. if I have millions of samples

## ▶ New Algorithm

\* DNN과 함께 사용!(Deep Neural Network)

여기엔 보통 엄청 많은 Data를 사용하기 때문에

Randomly choose an initial solution,  $w_0^0 w_1^0$

$t = 0$

New algorithm of GDM (sample 97일차)

Repeat

$$w_0^{t+1} = w_0^t - \eta \sum_{(\mathbf{x}_i, t_i) \in Data} \left. \frac{\partial E_i}{\partial w_0} \right|_{w_0=w_0^t, w_1=w_1^t}$$

*python에서 for, while로 구현*

$$w_1^{t+1} = w_1^t - \eta \sum_{(\mathbf{x}_i, t_i) \in Data} \left. \frac{\partial E_i}{\partial w_1} \right|_{w_0=w_0^t, w_1=w_1^t}$$

$t = t + 1$

Until stopping condition is satisfied

# But.. if I have millions of samples

## ▶ New Algorithm

Randomly choose an initial solution,  $w_0^0 w_1^0$

$t = 0$

Repeat

$$g_0^t = 0; g_1^t = 0$$

for all  $(x_i, t_i) \in Data$

$$g_0^t = g_0^t + \frac{\partial E_i}{\partial w_0} \Big|_{w_0=w_0^t, w_1=w_1^t}$$

$$g_1^t = g_1^t + \frac{\partial E_i}{\partial w_1} \Big|_{w_0=w_0^t, w_1=w_1^t}$$

$$w_0^{t+1} = w_0^t - \eta g_0^t$$

$$w_1^{t+1} = w_1^t - \eta g_1^t$$

$$t = t + 1$$

Until stopping condition is satisfied

]

$\sum \frac{\partial E_i}{\partial W}$  를 계산

# But.. if I have millions of samples

---

## ▶ Example

- ▶ Given data:  $\{(0.0, 0.0), (1.0, 1.0), (1.0, 2.0)\}$
- ▶ Model chosen:  $f(x; w_0, w_1) = w_1 x + w_0$
- ▶ Find  $w$  to minimize

$$E(w_0, w_1) = \sum_{(\mathbf{x}_i, t_i) \in Data} (t_i - f(\mathbf{x}_i; w_1, w_0))^2$$

# But.. if I have millions of samples

## ▶ Example

- ▶ Step 1: Formulate  $\frac{\partial E_i}{\partial w_0} \Big|_{w_0=w_0^t, w_1=w_1^t}$  and  $\frac{\partial E_i}{\partial w_1} \Big|_{w_0=w_0^t, w_1=w_1^t}$

$$g_0^t = g_0^t + \frac{\partial E_i}{\partial w_0} \Big|_{w_0=w_0^t, w_1=w_1^t} \quad g_1^t = g_1^t + \frac{\partial E_i}{\partial w_0} \Big|_{w_0=w_0^t, w_1=w_1^t}$$

$$E_i = (t_i - (w_1 x_i + w_0))^2$$

$$\frac{\partial E_i}{\partial w_0} = -2(t_i - (w_1 x_i + w_0))$$

$$\frac{\partial E_i}{\partial w_1} = -2(t_i - (w_1 x_i + w_0)) \cdot x_i$$

$$\frac{\partial E_i}{\partial w_0} \Big|_{w_0=\textcolor{red}{w}_0^t, w_1=\textcolor{blue}{w}_1^t} = -2(t_i - (\textcolor{blue}{w}_1^t x_i + \textcolor{red}{w}_0^t))$$

$$\frac{\partial E_i}{\partial w_1} \Big|_{w_0=\textcolor{red}{w}_0^t, w_1=\textcolor{blue}{w}_1^t} = -2(t_i - (\textcolor{blue}{w}_1^t x_i + \textcolor{red}{w}_0^t)) \cdot x_i$$

# But.. if I have millions of samples

## ▶ Example

### ▶ Step 2: Plug the formulae into the algorithm

Randomly choose an initial solution,  $w_0^0 w_1^0$

$t = 0$

Repeat

$$g_0^t = 0; g_1^t = 0$$

for all  $(x_i, t_i) \in Data$

$$g_0^t = g_0^t - 2(t_i - (w_1^t x_i + w_0^t))$$

$$g_1^t = g_1^t - 2(t_i - (w_1^t x_i + w_0^t)) \cdot x_i$$

$$w_0^{t+1} = w_0^t - \eta g_0^t$$

$$w_1^{t+1} = w_1^t - \eta g_1^t$$

$$t = t + 1$$

Until stopping condition is satisfied

# But.. if I have millions of samples

## ▶ Example

- ▶ Step 3: Randomly initialize  $w_0^t$  and  $w_1^t$ :  $w_0^0 = 1, w_1^0 = 1$   
Set  $\eta$  with a small value:  $\eta = 0.1$   
Run the algorithm

Randomly choose an initial solution,  $w_0^0 w_1^0$

$t = 0$

Repeat

$$g_0^t = 0; g_1^t = 0$$

for all  $(x_i, t_i) \in Data$

$$g_0^t = g_0^t - 2(t_i - (w_1^t x_i + w_0^t))$$

$$g_1^t = g_1^t - 2(t_i - (w_1^t x_i + w_0^t)) \cdot x_i$$

$$w_0^{t+1} = w_0^t - \eta g_0^t$$

$$w_1^{t+1} = w_1^t - \eta g_1^t$$

$$t = t + 1$$

Until stopping condition is satisfied

$$Data = \{(0.0, 0.0), (1.0, 1.0), (1.0, 2.0)\}$$

$$w_0^0 = 1, w_1^0 = 1$$

$$\begin{aligned} g_0^0 &= \sum - 2(t_i - (w_1^0 x_i + w_0^0)) \\ &= -2(0 - (0 + 1)) - 2(1 - (1 + 1)) \\ &\quad - 2(2 - (1 + 1)) = 4 \end{aligned}$$

$$\begin{aligned} g_1^0 &= \sum - 2(t_i - (w_1^0 x_i + w_0^0)) \cdot x_i \\ &= -2(0 - (0 + 1)) \cdot 0 - 2(1 - (1 + 1)) \cdot 1 \\ &\quad - 2(2 - (1 + 1)) \cdot 1 = 2 \end{aligned}$$

$$w_0^1 = 1 - 0.1 \cdot 4 = 0.6$$

$$w_1^1 = 1 - 0.1 \cdot 2 = 0.8$$



# Question and Answer