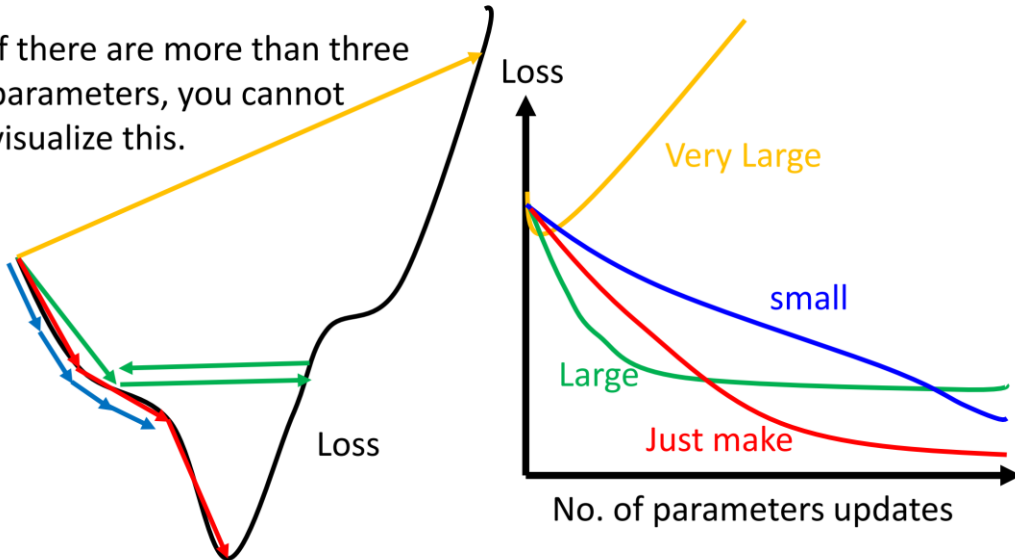


## Learning Rate

$$\theta^i = \theta^{i-1} - \eta \nabla L(\theta^{i-1})$$

Set the learning rate  $\eta$  carefully

If there are more than three parameters, you cannot visualize this.



But you can always visualize this.

## Adagrad

$$\eta^t = \frac{\eta}{\sqrt{t+1}} \quad g^t = \frac{\partial L(\theta^t)}{\partial w}$$

- Divide the learning rate of each parameter by the **root mean square of its previous derivatives**

### Vanilla Gradient descent

$$w^{t+1} \leftarrow w^t - \eta^t g^t$$

w is one parameters

### Adagrad

$$w^{t+1} \leftarrow w^t - \frac{\eta^t}{\sigma^t} g^t$$

$\sigma^t$ : **root mean square** of the previous derivatives of parameter w

Parameter dependent

## Adagrad

$\sigma^t$ : **root mean square** of the previous derivatives of parameter  $w$

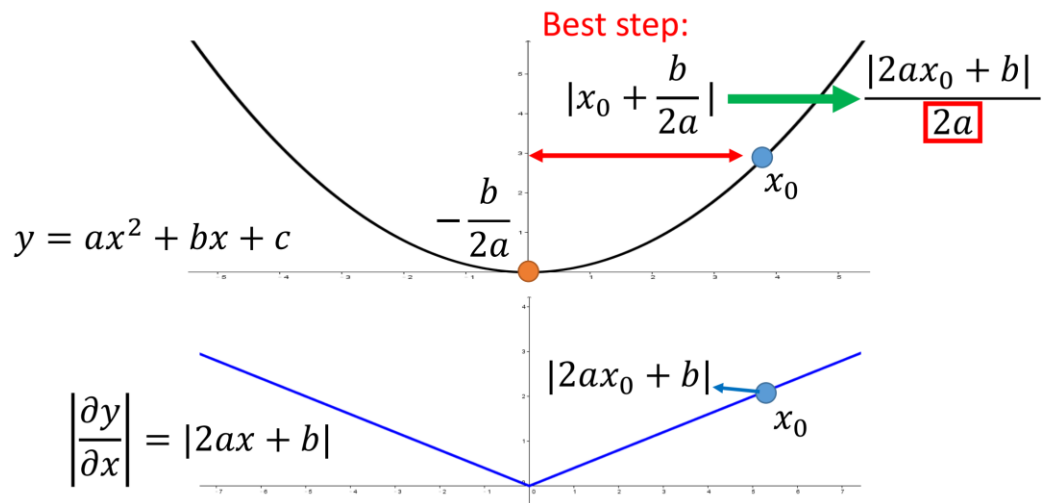
$$\begin{aligned}
 w^1 &\leftarrow w^0 - \frac{\eta^0}{\sigma^0} g^0 & \sigma^0 &= \sqrt{(g^0)^2} \\
 w^2 &\leftarrow w^1 - \frac{\eta^1}{\sigma^1} g^1 & \sigma^1 &= \sqrt{\frac{1}{2} [(g^0)^2 + (g^1)^2]} \\
 w^3 &\leftarrow w^2 - \frac{\eta^2}{\sigma^2} g^2 & \sigma^2 &= \sqrt{\frac{1}{3} [(g^0)^2 + (g^1)^2 + (g^2)^2]} \\
 &\vdots & & \\
 w^{t+1} &\leftarrow w^t - \frac{\eta^t}{\sigma^t} g^t & \sigma^t &= \sqrt{\frac{1}{t+1} \sum_{i=0}^t (g^i)^2}
 \end{aligned}$$

## Adagrad

- Divide the learning rate of each parameter by the **root mean square of its previous derivatives**

$$\begin{aligned}
 w^{t+1} &\leftarrow w^t - \frac{\eta^t}{\sigma^t} g^t \\
 &\quad \eta^t = \frac{\eta}{\sqrt{t+1}} \quad \text{1/t decay} \\
 &\quad \sigma^t = \sqrt{\frac{1}{t+1} \sum_{i=0}^t (g^i)^2} \\
 \Downarrow \\
 w^{t+1} &\leftarrow w^t - \frac{\eta}{\sqrt{\sum_{i=0}^t (g^i)^2}} g^t
 \end{aligned}$$

# Second Derivative



$$\frac{\partial^2 y}{\partial x^2} = 2a$$

The best step is

$\frac{|\text{First derivative}|}{\text{Second derivative}}$