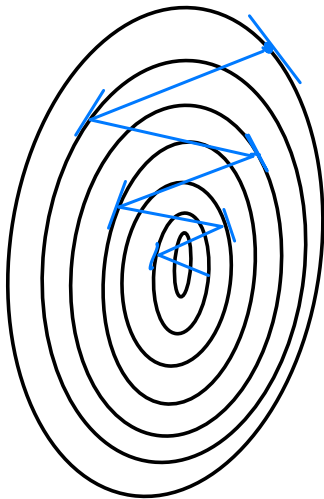


# #optimizer.

## 1) Momentum optimizer.



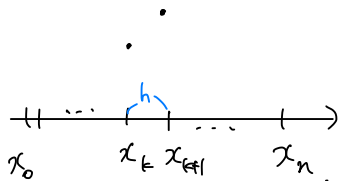
avoid zigzag.

$$X_{k+1} = X_k - s \cdot \underbrace{Z_k}_{\text{new direction.}}$$

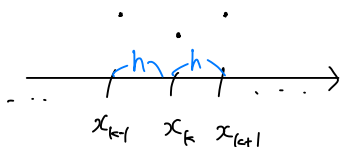
(if  $Z_k = \nabla f$  : simple gradient descent)

$$Z_k = \nabla f + \beta \cdot \underbrace{Z_{k-1}}_{\text{memory of previous step (momentum).}}$$

check.



$$f'(x_k)? \Rightarrow f'(x_k) = \frac{f(x_{k+1}) - f(x_k)}{h}$$



$$\begin{aligned} f''(x_k)? &\Rightarrow f''(x_k) = \frac{f'(x_k) - f'(x_{k-1}))}{h} \\ &= \frac{1}{h} \left\{ \frac{f(x_{k+1}) - f(x_k)}{h} - \frac{f(x_k) - f(x_{k-1}))}{h} \right\} \\ &= \frac{f(x_{k+1}) - 2f(x_k) + f(x_{k-1}))}{h^2} \end{aligned}$$

$$\begin{cases} \underline{X_{k+1}} = \underline{X_k} - s \cdot \nabla f & \longrightarrow 2 \text{ step} \approx \text{first order differential equation.} \\ \underline{X_{k+1}} = \underline{X_k} - s (\nabla f + \beta \cdot \underline{Z_k}) & \longrightarrow 3 \text{ step} \approx \text{second order differential equation.} \end{cases}$$

( $k \neq$  time variable but analogy)

$$\left. \begin{aligned} X_{k+1} &= X_k - s \cdot Z_k \\ Z_k &= \nabla f_k + \beta \cdot Z_{k-1} \end{aligned} \right) \longrightarrow \left. \begin{aligned} X_{k+1} &= X_k - s \cdot Z_k \\ Z_{k+1} - \nabla f_k &= \beta \cdot Z_k \end{aligned} \right) \longrightarrow \left. \begin{aligned} X_{k+1} &= X_k - s \cdot Z_k \\ Z_{k+1} - s X_{k+1} &= \beta \cdot Z_k \end{aligned} \right)$$

$k \rightarrow k+1$ .

$H_f$ .

$$f = \frac{1}{2} \cdot x^T S x$$

$$\nabla f = Sx$$

$$\Rightarrow \begin{bmatrix} 1 & 0 \\ -s & 1 \end{bmatrix} \begin{bmatrix} X \\ Z \end{bmatrix}_{k+1} = \begin{bmatrix} 1 & \text{step size} \\ 0 & \beta \end{bmatrix} \begin{bmatrix} X \\ Z \end{bmatrix}_k$$

$$\left( \begin{aligned} Sx &= \lambda x \rightarrow x_k = c_k \cdot x \\ Z_k &= d_k \cdot x \\ Sx_k &= c_k \cdot \lambda \cdot x \end{aligned} \right)$$

$$\Rightarrow \begin{bmatrix} 1 & 0 \\ -s & 1 \end{bmatrix} \begin{bmatrix} x \\ z \end{bmatrix}_{k+1} = \begin{bmatrix} 1 & -s \\ 0 & \beta \end{bmatrix} \begin{bmatrix} x \\ z \end{bmatrix}_k$$

$$\Rightarrow \begin{bmatrix} 1 & 0 \\ -\lambda & 1 \end{bmatrix} \begin{bmatrix} c_k \\ d_{k+1} \end{bmatrix} = \begin{bmatrix} 1 & -s \\ 0 & \beta \end{bmatrix} \begin{bmatrix} c_k \\ d_k \end{bmatrix}$$

choose  $s, \beta$ .

$$\begin{aligned} \begin{bmatrix} c_{k+1} \\ d_{k+1} \end{bmatrix} &= \begin{bmatrix} 1 & 0 \\ \lambda & 1 \end{bmatrix} \begin{bmatrix} 1 & -s \\ 0 & \beta \end{bmatrix} \begin{bmatrix} c_k \\ d_k \end{bmatrix} \\ &= \begin{bmatrix} 1 & -s \\ \lambda & \beta - \lambda s \end{bmatrix} \begin{bmatrix} c_k \\ d_k \end{bmatrix} \\ &= \underline{R} \cdot \begin{bmatrix} c_k \\ d_k \end{bmatrix} \end{aligned}$$

$\lambda$ : Sol e.value.

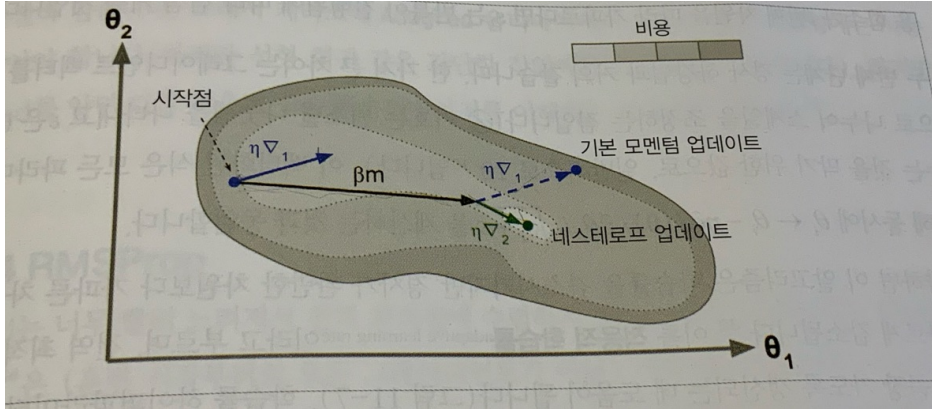
$\begin{bmatrix} \lambda_{\min} \\ \lambda_{\max} \end{bmatrix}$  = condition number.

$$\underline{X}_{k+1} = \underline{R} \cdot \underline{X}_k$$

we want to choose  $s, \beta$

to make e.values of  $R$  as small as possible.

2) Nesterov momentum optimization.



$$X_{k+1} = X_k + \beta(x_k - x_{k-1}) - s \cdot \nabla f(x_k + t(x_k - x_{k-1}))$$

$$\Rightarrow \begin{bmatrix} C_{KH} \\ d_{KH} \end{bmatrix} = \begin{bmatrix} 0 & (1 - \zeta \cdot \lambda) \\ -\beta & (1 + \beta)(1 - \zeta \cdot \lambda) \end{bmatrix} \begin{bmatrix} C_K \\ d_K \end{bmatrix}$$

$$= R \cdot \begin{bmatrix} C_K \\ d_K \end{bmatrix}$$

3) Adaptive Method.

$$x_{k+1} = x_k - \underset{\uparrow}{s} \nabla f.$$

$$x_{k+1} = x_k - s_k \cdot D_k.$$

$$\begin{cases} D_k = D(\nabla L_k, \nabla L_{k-1}, \dots, \nabla L_0) \\ s_k = s(\nabla L_k, \nabla L_{k-1}, \dots, \nabla L_0) \end{cases} \Rightarrow D_k = \nabla L_k.$$

i) AdaGrad.

$$s \leftarrow s + \nabla_{\theta} J(\theta) \otimes \nabla_{\theta} J(\theta) \longrightarrow \sum_{i=1}^k (\nabla_{\theta} J(\theta) \otimes \nabla_{\theta} J(\theta))$$

$$\theta \leftarrow \theta - \frac{\nabla_{\theta} J(\theta)}{\sqrt{s}} \otimes \sqrt{s} \epsilon.$$

①  $\rightarrow$   $\frac{\nabla_{\theta} J(\theta)}{\sqrt{s}}$   $\rightarrow$   $\frac{\nabla_{\theta} J(\theta)}{\sqrt{s}}$   $\rightarrow$   $\frac{\nabla_{\theta} J(\theta)}{\sqrt{s}}$   $\rightarrow$   $\frac{\nabla_{\theta} J(\theta)}{\sqrt{s}}$

ii) RMSprop.

exponential decay.

$$s \leftarrow \beta s + (1-\beta) \nabla_{\theta} J(\theta) \otimes \nabla_{\theta} J(\theta) \longrightarrow (1-\beta) \sum_{i=1}^k \beta^{k-i} (\nabla_{\theta} J(\theta) \otimes \nabla_{\theta} J(\theta))$$

$$\theta \leftarrow \theta - \frac{\nabla_{\theta} J(\theta)}{\sqrt{s}} \otimes \sqrt{s} \epsilon.$$

# Adam. (Adaptive + Momentum)

$$\downarrow m \leftarrow \beta_1 m - \underbrace{(1-\beta_1)}_{0.1} \nabla_{\theta} J(\theta).$$

$$\downarrow S \leftarrow \beta_2 S + (1-\beta_2) \nabla_{\theta} J(\theta) \otimes \nabla_{\theta} J(\theta).$$

$$\hat{m} \leftarrow \frac{m}{1-\beta_1^t}$$

$$\hat{S} \leftarrow \frac{S}{1-\beta_2^t}$$

$$\theta \leftarrow \theta + \eta \hat{m} \oslash \sqrt{\hat{S} + \epsilon}.$$

---

$$\beta_1 = 0.9$$

$$m_0 = 0$$

## # Regularization

-  $l_1, l_2$

- drop out

- max\_norm

