



Deep Learning Basic

Jaewon Kim, Dankook Univ.

Chapter 4-2



Contents

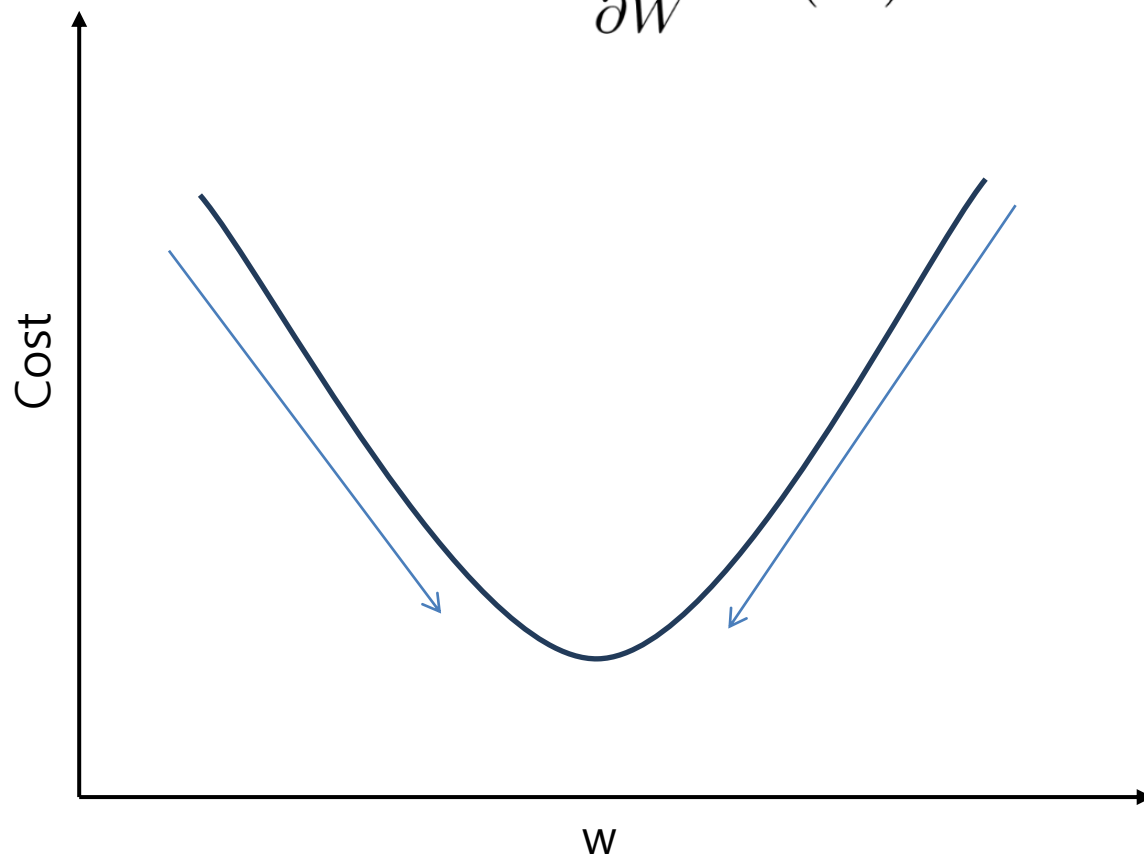
Part 1. Optimizer

- Gradient Descent Review
- Kinds of Optimizer



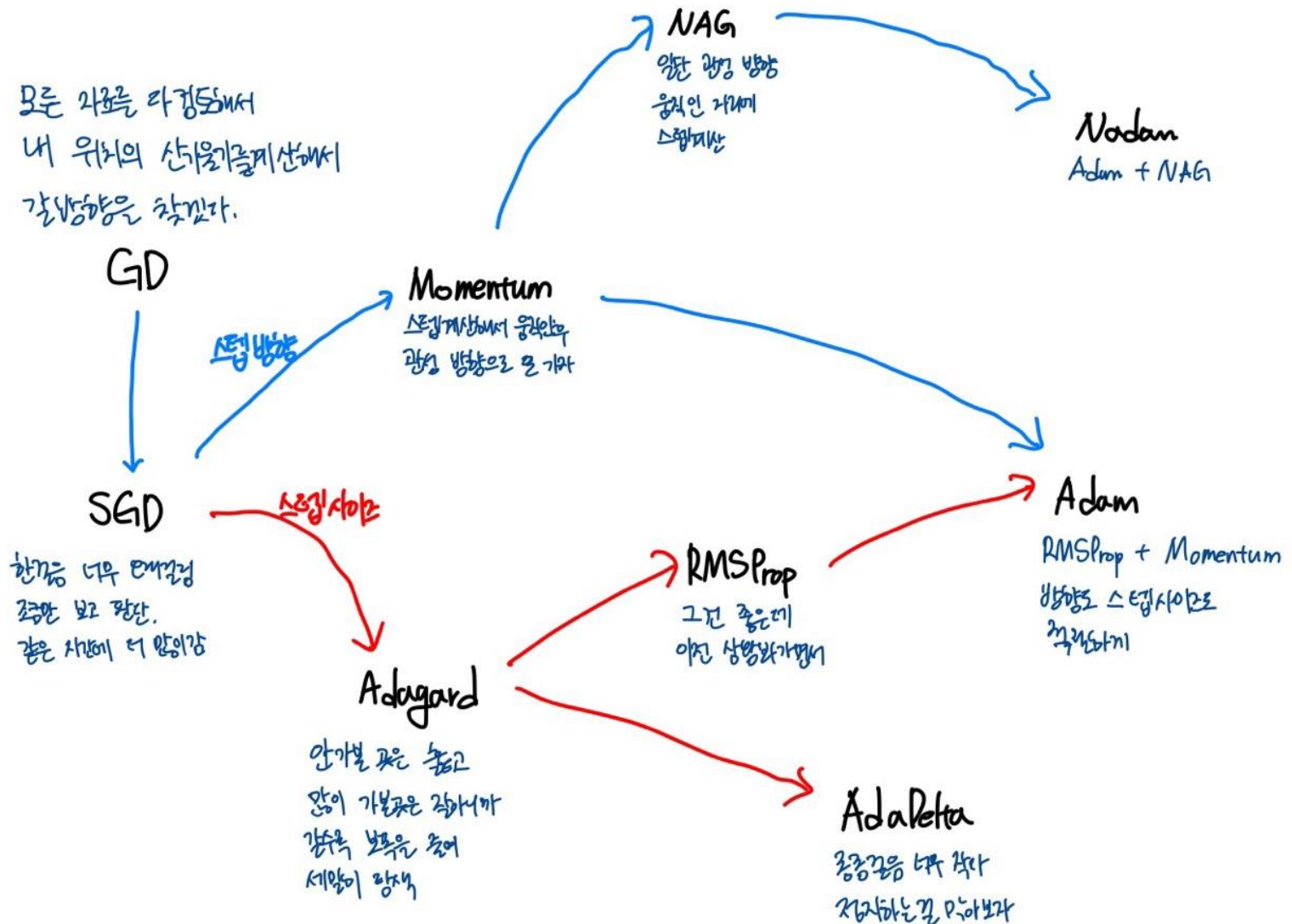
Optimization (Gradient Decent)

$$W := W - \alpha \frac{\partial}{\partial W} cost(W)$$

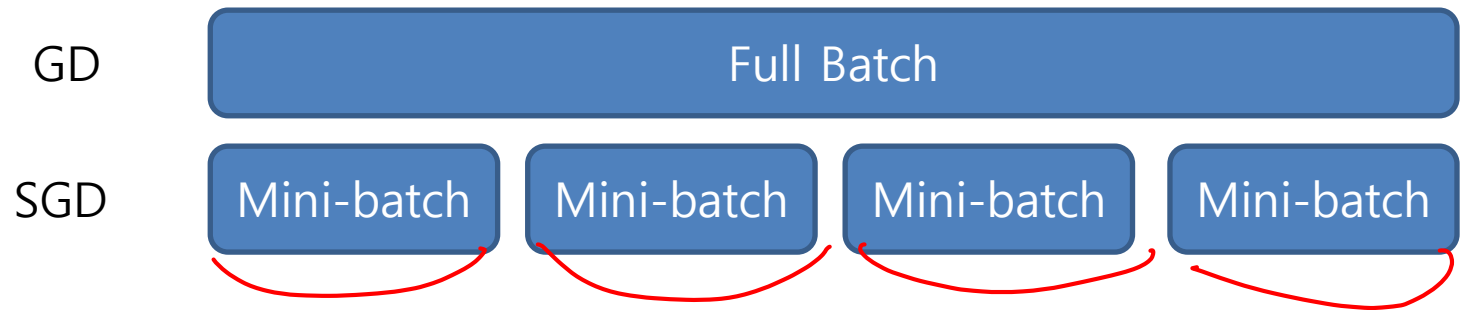


But too slow, local minimum

Kinds of Optimizer



Stochastic Gradient Decent (SGD)



batch size 2 이상



local minimum에 빠지기 쉬움.

Momentum

- SGD에 Momentum 개념을 추가
SGD

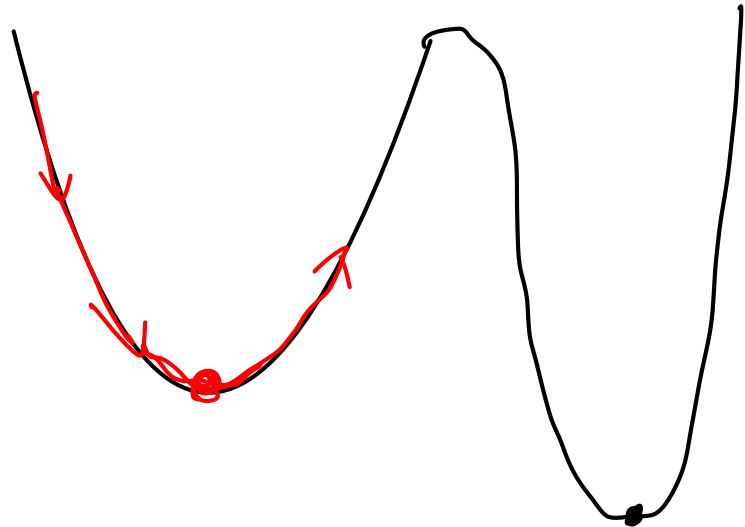
$$v \leftarrow \alpha v - \eta \frac{\partial L}{\partial W}$$

이전의 값. 가속도계 (0.9, 1 이하의 값)

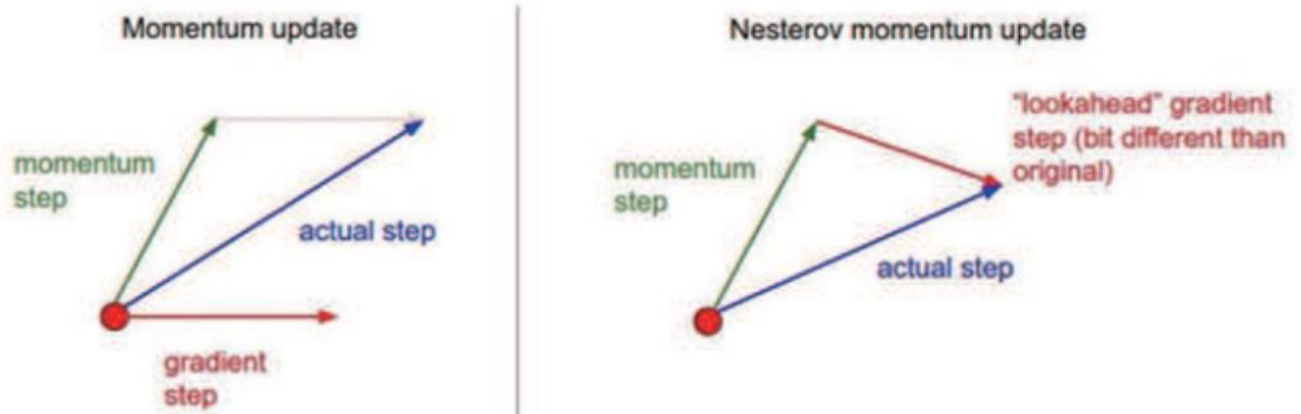
$$W \leftarrow W + v$$

이전 batch x , 이전의 계수로 반영.

안정성



NAG (Nesterov Accelerated Gradient)



$$\theta = \theta - v_t$$

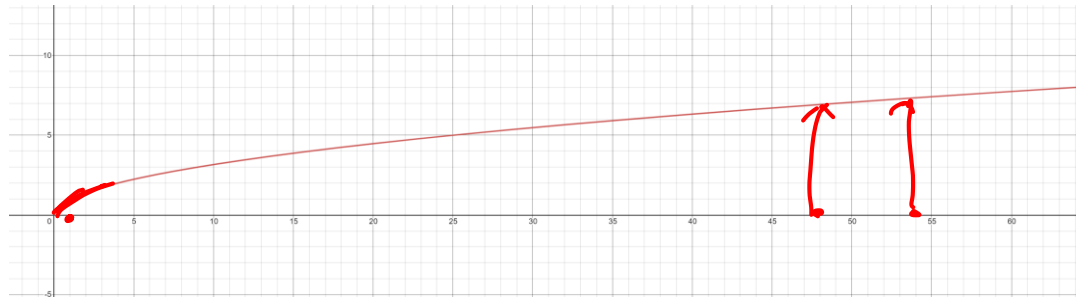
$$v_t = \gamma v_{t-1} + \eta \nabla_{\theta} J(\theta - \gamma v_{t-1})$$

Adagrad (Adaptive Gradient)

$$\theta_{t+1} = \theta - \frac{\eta}{\sqrt{G_t + \epsilon}} \cdot \nabla_{\theta} J(\theta_t)$$

$$G_t = G_{t-1} + (\nabla_{\theta} J(\theta_t))^2$$

→ 얼마나 많이 변했는가를.



$$y = \sqrt{x}$$

• 크게 변하면 η 감소
• 작게 변하면 η 증가

RMSProp

$$h_i \leftarrow \rho h_{i-1} + (1 - \rho) \frac{\partial L_i}{\partial W} \odot \frac{\partial L_i}{\partial W}$$

~~decay constant~~ (0 ~ 1) 사이의 값

1. 초기값을 0으로 설정.
2. ρ hyperparameter 지정.

AdaDelta

$$\begin{aligned}\theta_{t+1} &= \theta_t - \Delta\theta \\ \Delta\theta &= \frac{\sqrt{s + \epsilon}}{\sqrt{G + \epsilon}} \cdot \nabla_{\theta} J(\theta_t) \\ s_{t+1} &= \gamma s_t + (1 - \gamma) \Delta\theta^2 \\ G_{t+1} &= \gamma G_t + (1 - \gamma) (\nabla_{\theta} J(\theta_t))^2\end{aligned}$$

파라미터가

오염이나 변하는지

관리,

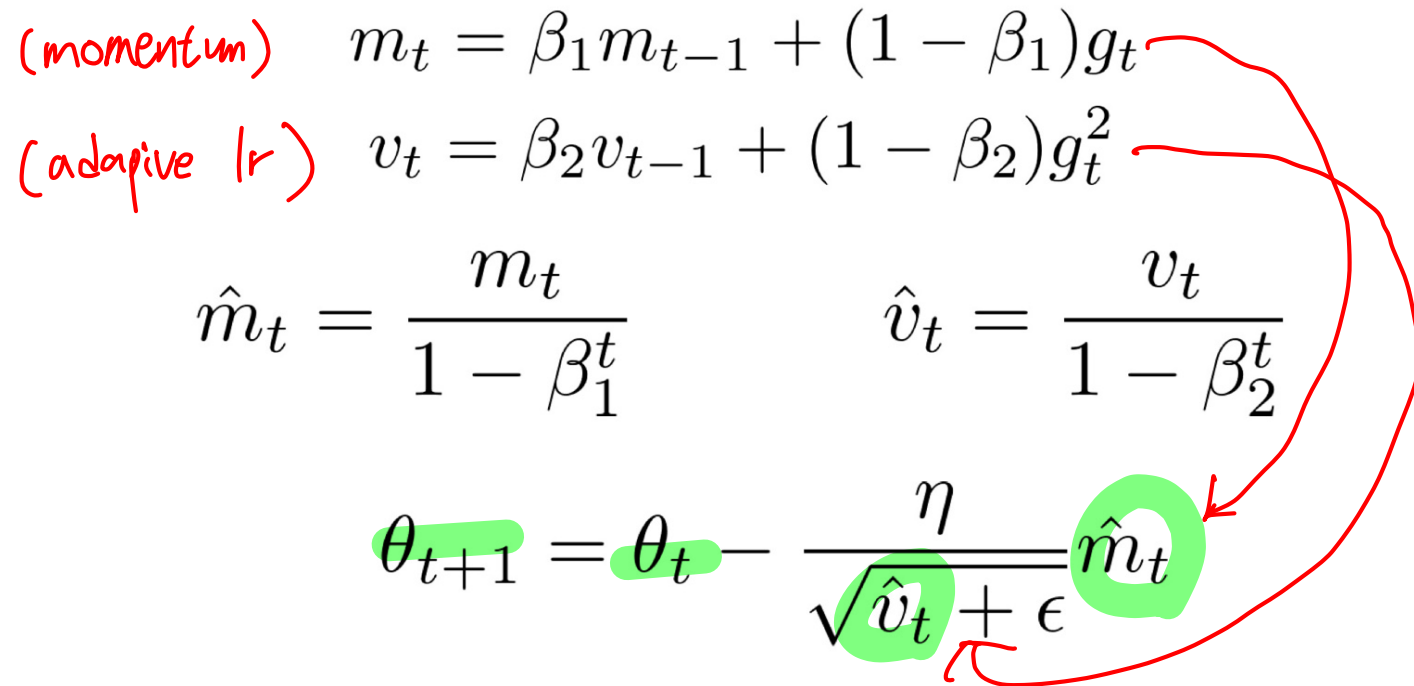
Adam

(momentum) $m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$

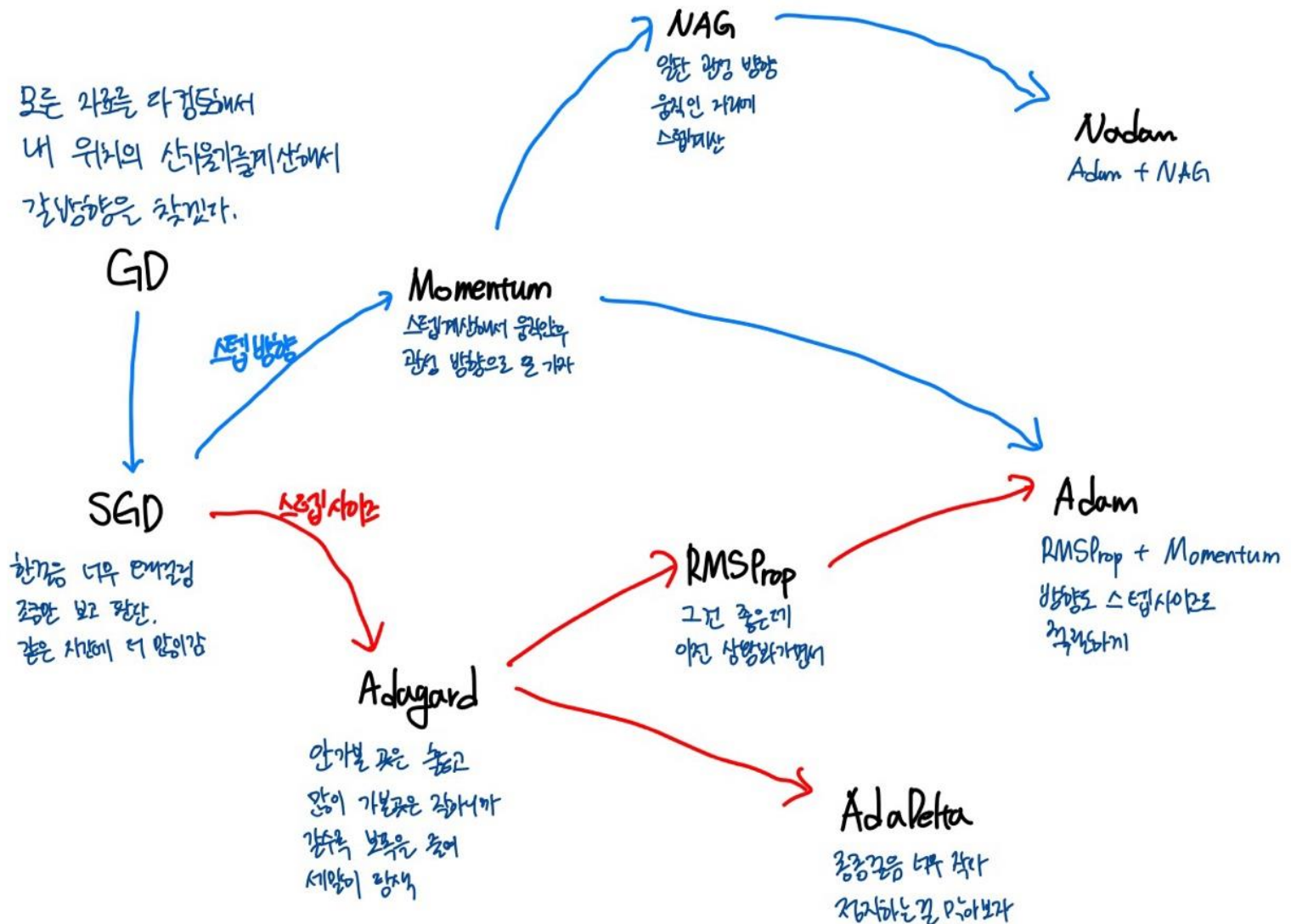
(adaptive lr) $v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}$$

$$\hat{v}_t = \frac{v_t}{1 - \beta_2^t}$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \hat{m}_t$$


Kinds of Optimizer



Optimizer

<https://gomguard.tistory.com/187>

Thank you...!!!