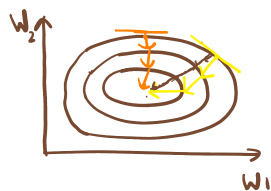


# Adaptive Gradient (Adagrad)

Saturday, October 9, 2021 2:18 PM



momentum

$$\eta = \beta m + \eta \frac{\partial L}{\partial W} \Big|_{W=W_0} \rightarrow \textcircled{m}$$

NAG

$$\eta = \beta m + \eta \frac{\partial L}{\partial W} \Big|_{W=(W_0 - \beta m)} \rightarrow \textcircled{NAG}$$

$\therefore m$  in both cases is proportional to gradient

here direction is not much of a concern for  $m$

$J(W_1, W_2)$



Adagrad  $\Rightarrow$  correct the direction at initial steps

Algorithm  $\Rightarrow \beta_t$  - scaling factor  $\Rightarrow \beta_t + \left(\frac{\partial L}{\partial W}\right)^2$  where initialize  $\beta = 0 \rightarrow \textcircled{1}$

$$\eta^* = \frac{\eta}{\sqrt{\beta_t + \epsilon}}$$

$\beta$  - scaling factor

$\epsilon$  - Smoothing factor (avoid division by 0)

$$W_t = W_0 - \eta \frac{\frac{\partial L}{\partial W} \Big|_W}{\sqrt{\beta_t + \epsilon}}$$

$\rightarrow \textcircled{11}$

$$\frac{\eta}{\sqrt{\beta_t + \epsilon}} \cdot \frac{\partial L}{\partial W}$$

Learning rate

Q - Why we need the scaling factor??

$$\beta_t = \beta_{t-1} + \underbrace{\nabla_0 J(\theta) * \nabla_0 J(\theta)}_{\text{square term - magnitude of vectors}} \rightarrow \textcircled{1}$$

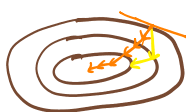
$$\frac{\partial L}{\partial W} = \nabla_0 J(\theta)$$

$$W = W - \eta \frac{\nabla_0 J(\theta)}{\sqrt{\beta_t + \epsilon}} \rightarrow \text{vector} \rightarrow \textcircled{1}$$

$$\text{Unit vector} = \frac{\text{Vector}}{|\text{Vector}|}$$

$$W \Rightarrow W - \eta \frac{\nabla_0 J(\theta)}{\sqrt{\beta_t}} \Rightarrow W - \eta \frac{\nabla_0 J(\theta)}{\sqrt{|\nabla_0 J(\theta)|^2}} = W - \eta \frac{\nabla_0 J(\theta)}{|\nabla_0 J(\theta)|}$$

Unit vector  
(Scaling to unit)



$$\text{Step 1} \Rightarrow W - \eta |\text{unit vector}|$$

Small value (smaller than  $\frac{\partial L}{\partial W}$ )  
 $\frac{\partial L}{\partial W}^2$

Step 1  $\Rightarrow \beta_0 = 0, W = W_0$

$$\beta_1 = \beta_0 + \left(\frac{\partial L}{\partial W}\right)^2 \Rightarrow \beta_1 = \left(\frac{\partial L}{\partial W} \Big|_{W=W_0}\right)^2$$

$$W_1 = W_0 - \eta \frac{\frac{\partial L}{\partial W}}{\sqrt{\beta_1 + \epsilon}} = W_0 - \eta \frac{\frac{\partial L}{\partial W} \Big|_{W=W_0}}{\sqrt{\left(\frac{\partial L}{\partial W} \Big|_{W=W_0}\right)^2 + \epsilon}} \rightarrow \textcircled{10}$$

$$\text{Step 2} \Rightarrow \beta_1 = \left(\frac{\partial L}{\partial W} \Big|_{W=W_0}\right)^2, W = W_1$$

$$\beta_2 = \beta_1 + \left(\frac{\partial L}{\partial W} \Big|_{W=W_1}\right)^2 \Rightarrow \underbrace{\left(\frac{\partial L}{\partial W} \Big|_{W=W_0}\right)^2}_{\text{Previous}} + \underbrace{\left(\frac{\partial L}{\partial W} \Big|_{W=W_1}\right)^2}_{\text{Current}}$$

$$W_2 = W_1 - \eta \frac{\frac{\partial L}{\partial W} \Big|_{W=W_1}}{\sqrt{\beta_2 + \epsilon}} \rightarrow \textcircled{11}$$

$$\Rightarrow \sqrt{\beta_2 + \epsilon} > \sqrt{\beta_1 + \epsilon}$$

gradient is adapting based on loss function

Similarly we can generalize (11)

#### Advantage of Adagrad

- 1) Initially adjust the direction
- 2) Less tuning of learning rate

#### Disadvantages

- 1) Stops early before reaching global minimum due to decaying in gradient
- 2) Takes longer time to converge \* Not recommended \*