

# Rms Prop - Root mean square propagation

Saturday, October 9, 2021 3:43 PM

\* Solves the problem of early stopping of Adagrad by accumulating gradients from recent iteration by using exponential decay

$$\beta_t = \beta_{t-1} + \left( \frac{\partial L}{\partial W} \Big|_{W=W_{t-1}} \right)^2 \rightarrow \textcircled{1} \text{ from Adagrad}$$

$$\beta_t = \gamma \beta_{t-1} + (1-\gamma) \left( \frac{\partial L}{\partial W} \right)^2 \rightarrow \textcircled{2} \text{ the change}$$

$$W = W - \eta \frac{\frac{\partial L}{\partial W}}{\sqrt{\beta_t + \epsilon}} \Rightarrow W_t = W_{t-1} - \eta \frac{\frac{\partial L}{\partial W} \Big|_{W=W_{t-1}}}{\sqrt{\gamma \beta_{t-1} + (1-\gamma) \left( \frac{\partial L}{\partial W} \right)^2}}$$

$$\beta = 0.9 \text{ works well}$$

Step 1  $\Rightarrow \beta_0 = 0, W = W_0, \gamma = 0.9$

$$\beta_1 = \gamma \beta_0 + (1-\gamma) \left( \frac{\partial L}{\partial W} \right)^2 \rightarrow \textcircled{1}$$

$$\beta_1 = \gamma \cdot 0 + (1-\gamma) \left[ \frac{\partial L}{\partial W} \Big|_{W_0} \right]^2$$

$$W_1 = W_0 - \eta \frac{\frac{\partial L}{\partial W} \Big|_{W_0}}{\sqrt{(1-\gamma) \left[ \frac{\partial L}{\partial W} \Big|_{W_0} \right]^2 + \epsilon}} \rightarrow \textcircled{1}$$

Step 2  $\Rightarrow \beta_1 = 0, W = W_1, \gamma = 0.9$

$$\beta_2 = \gamma \beta_1 + (1-\gamma) \left( \frac{\partial L}{\partial W} \Big|_{W_1} \right)^2$$

$$\beta_2 = \gamma \left[ (1-\gamma) \left[ \frac{\partial L}{\partial W} \Big|_{W_0} \right]^2 \right] + (1-\gamma) \left[ \frac{\partial L}{\partial W} \Big|_{W_1} \right]^2$$

$$(1-\gamma) \left[ \underbrace{\gamma}_{\substack{\text{90\% of past} \\ \text{0.9}}} \left[ \frac{\partial L}{\partial W} \Big|_{W_0} \right]^2 + \underbrace{\left[ \frac{\partial L}{\partial W} \Big|_{W_1} \right]^2}_{\substack{\text{100\% of current} \\ \text{1.0}}} \right]$$

Step 1  $\Rightarrow \beta_2 \rightarrow \odot$ ,  $w = w_2$ ,  $\gamma = 0.9$

$$\beta_3 = \gamma \beta_2 + (1-\gamma) \left( \frac{\partial L}{\partial w} \Big|_{w=w_2} \right)^2$$

$$\beta_3 = (1-\gamma) \left[ \gamma^2 \left( \frac{\partial L}{\partial w} \Big|_{w=w_0} \right)^2 + \gamma \left( \frac{\partial L}{\partial w} \Big|_{w=w_1} \right)^2 + \left( \frac{\partial L}{\partial w} \Big|_{w=w_2} \right)^2 \right]$$

$$w_3 = w_2 - \eta \frac{\frac{\partial L}{\partial w} \Big|_{w_2}}{\sqrt{\beta_3 + \epsilon}} \rightarrow \textcircled{11}$$

Adagrad

$$\beta_{2, \text{Ada}} = \left( \frac{\partial L}{\partial w} \Big|_{w=w_0} \right)^2 + \left( \frac{\partial L}{\partial w} \Big|_{w=w_1} \right)^2$$

Rmsprop

$$\beta_{2, \text{RMS}} = (1-\gamma) \left[ \gamma \left( \frac{\partial L}{\partial w} \Big|_{w=w_0} \right)^2 + \left( \frac{\partial L}{\partial w} \Big|_{w_1} \right)^2 \right]$$

$\beta_{2, \text{Ada}} \geq \beta_{2, \text{Rmsprop}} \Rightarrow$  But  $\beta_2$  is a denominator for weight update  
So Rms prop is little faster to perform early stopping