

SGD

The gradient updates proceed as follows in each case :-

For all cases, batch size = m ,
SGD :- $\odot \rightarrow$ Hadamard product

$$\text{Gradient :- } \nabla_{\theta}(R(\theta)) = \nabla_{\theta} \left(\sum_{i=1}^m R_i(\theta) \right)$$

$$\text{update :- } \theta^{(r+1)} = \theta^{(r)} - (\eta) \nabla_{\theta}(R(\theta))$$

Momentum :-

$$v^{(0)} = 0$$

$$v^{(r)} = \alpha v^{(r-1)} - \sum \nabla_{\theta}^{(r)} R(\theta)$$

$$\theta^{(r)} = \theta^{(r-1)} + v^{(r)}$$

$v \rightarrow$ similar to velocity.

Clearly, ~~the~~ effective learning rate depends on the parameters in previous step.

Nesterov :-

> Here, we apply an interim update on the parameter set and find gradient wrt that update.

$$\begin{aligned}\theta_{(r)} &= \theta_{(r)} + \alpha v \\ g &= \left(\frac{1}{m}\right) \nabla_{\theta} \sum R_i(\theta) \\ v &= \alpha v - \epsilon g \\ \theta_{(r+1)} &= \theta_{(r)} + v_{(r)}\end{aligned}$$

Ada-grad :-

Find the gradient

$$g = \frac{1}{m} \nabla_{\theta} \sum R_i(\theta)$$

$$r = r + \underbrace{g \odot g}_{\text{square of gradient.}}$$

update :-

$$\theta = \theta - \frac{\epsilon}{\delta + \sqrt{r}} \cdot g$$

$$S = 0 \quad (1e-10)$$

RHS - Prop:

Small variation of Ada-grad.

$$r^{(k+1)} = r^{(k)} \cdot (1-p)$$

$$g = \frac{1}{m} \nabla_{\theta} \sum R_i(\theta)$$

$$r^{(k+1)} = r^{(k)} (1-p) + p \cdot g \odot g$$

update :-

$$\theta^{(k+1)} = \theta^{(k)} - \frac{\epsilon}{\sqrt{S + r^{(k+1)}}} \odot g$$

Adamm :-

We decay the learning rate over batches using the 1st and 2nd moments.

So:-

$$g = \frac{1}{m} \nabla_{\theta} (\sum R_i(\theta))$$

$t = t+1 \longrightarrow$ time step.

$$S = p_1 S + (1-p_1) g \odot g \quad (1^{st} \text{ moment})$$

$$r = p_2 r + (1-p_2) g \odot g \quad (2^{nd} \text{ "})$$

update :-

$$\theta = \theta - \varepsilon \frac{\hat{s}}{\sqrt{\hat{r} + \delta}}$$

where

$$\hat{s} = \frac{s}{1 - \rho_1^t}$$

$$\hat{r} = \frac{r}{1 - \rho_2^t}$$