

Adam weight update rule

We take $v(0)$ and $s(0)$ as zero

At iteration t

Compute $\partial C / \partial w$ on a mini-batch

$$v(t) := \beta_1 v(t) + (1 - \beta_1) \partial C / \partial w$$

$$s(t) := \beta_2 s(t) + (1 - \beta_2) \frac{\delta C^2}{\delta w}$$

$$v_t^{\text{corrected}} = \frac{v_t}{(1 - \beta_1)}$$

$$s_t^{\text{corrected}} = \frac{s_t}{(1 - \beta_2)}$$

$$w := w - \alpha \frac{v_t^{\text{corrected}}}{\sqrt{s_t^{\text{corrected}} + \epsilon}}$$

$\beta_1 \in (0, 1)$ Generally, it is taken to be 0.9

$\beta_2 \in (0, 1)$ Generally, it is taken to be 0.99

ϵ is taken to be 10^{-8}

Here α is learning rate.

We tune α to get optimum results.