

Gradient Descent

$$\theta := \theta - \alpha d\theta$$

hyperparameter: α

Gradient Descent with Momentum

$$v_\theta := \beta v_\theta + (1 - \beta) d\theta$$

hyperparameter: α, β

$$\theta := \theta - \alpha v_\theta$$

common value: 0.9

some literature

$$\left(v_\theta := \beta v_\theta + d\theta \right)$$

RMSprop (Root Mean Square)

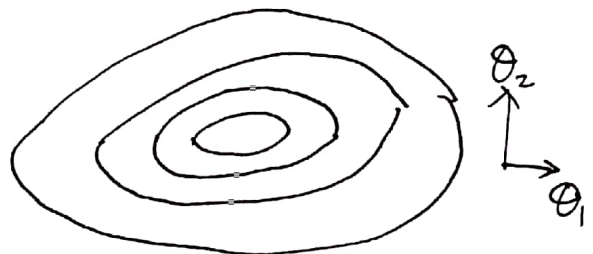
$$s_\theta := \beta s_\theta + (1 - \beta) d\theta^2$$

element-wise

$$\theta := \theta - \alpha \frac{d\theta}{\sqrt{s_\theta + \epsilon}}$$

ensure numerical stability

intuition:



in θ_1 dimension,

$d\theta_1$ smaller $\rightarrow s_{\theta_1}$ smaller $\rightarrow \frac{1}{\sqrt{s_{\theta_1}}}$ larger

in θ_2 dimension

$d\theta_2$ larger $\rightarrow s_{\theta_2}$ larger $\rightarrow \frac{1}{\sqrt{s_{\theta_2}}}$ smaller

Adam Optimization (Adaptive Moment Estimation)

(gradient descent with momentum
+
RMS prop gradient descent)

$$V_{\theta} := \beta_1 V_{\theta} + (1 - \beta_1) d\theta$$

$$S_{\theta} := \beta_2 S_{\theta} + (1 - \beta_2) d\theta^2$$

$$V_{\theta \text{ correct}} = \frac{V_{\theta}}{1 - \beta_1^t}$$

$$S_{\theta \text{ correct}} = \frac{S_{\theta}}{1 - \beta_2^t}$$

$$\theta := \theta - \alpha \frac{V_{\theta \text{ correct}}}{\sqrt{S_{\theta \text{ correct}} + \epsilon}}$$

Hyperparameters:

α

β_1 0.9

β_2 0.999

ϵ : 10^{-8}

Learning Rate Decay

$$- \alpha = \frac{\alpha_0}{1 + (\text{decay-rate})^{\text{epoch-num}}}$$

$$- \alpha = 0.95^{\text{epoch-num}} * \alpha_0$$

$$- \alpha = \frac{k}{\sqrt{\text{epoch-num}}} \alpha_0$$

- manually adjust α

Exponentially Weighted Average

$$\hat{v}_t = \beta \hat{v}_{t-1} + (1-\beta) O_t$$

Bias Correction

$$\frac{\hat{v}_t}{1-\beta^t}$$

correct the
initial phase

