

Lecture 14. Advanced opt. Algos. II.

RMSprop

$$w^{i+1} = w^i + \underline{\Delta w^i}$$

$$\underline{\Delta w^i} = \frac{-\epsilon}{\sqrt{r^i + \delta}} (\nabla_{w^i} L)$$

$$r^i = \rho \cdot r^{i-1} + (1-\rho) (\nabla_{w^i} L)^2 \quad 0 < \rho < 1$$

δ : small positive constant.
 i is large $\rightarrow e^{i-1} \rightarrow 0$

$$r^i = \rho^{i-1} r^1 + \rho^{i-2} r^2$$

1. RMSprop with Nesterov momentum.

$$\underline{\Delta w^i} = 2 \cdot \Delta w^{i-1} + \epsilon / \sqrt{r^i} (-\nabla_{w^i} L)$$

$$r^i = \rho \cdot r^{i-1} + (1-\rho) (\nabla L)^2$$

2. Adaptive moments (Adam, Kingma and Ba, 2014)

ideal is to combine RMSprop + Momentum.

$$\Delta w^i = \frac{-\epsilon}{\sqrt{\hat{r}^i}} \cdot \hat{s}^i \quad 0 < \rho_1 < 1, 0 < \rho_2 < 1$$

$$r^i = \rho_2 \cdot r^{i-1} + (1 - \rho_2) \cdot (\nabla_w L)^2: \text{2nd moment variable.}$$

$$s^i = \rho_1 \cdot s^{i-1} + (1 - \rho_1) \cdot (\nabla_w L): \text{1st moment variable.}$$

$$\hat{r}^i = \frac{r^i}{1 - (\rho_2)^i} \rightarrow \text{power.} \quad \hat{s}^i = \frac{s^i}{1 - (\rho_1)^i}$$

$$\text{~~init~~ } r^0 = s^0 = 0$$

$$i \rightarrow \text{large} \rightarrow +\infty \rightarrow (\rho_2)^i = 0, (\rho_1)^i = 0$$

\hat{r}^i & \hat{s}^i introduced to correct bias
at the beginning of training with small initial w .
 $\hat{r}^i = r^i \quad \hat{s}^i = s^i$

3. $w^i = w^{i-1} + \Delta w^i \rightarrow$ combined step size

Algorithm	Δw^i	pros	Cons.
SGD	$\Delta w^i = -\epsilon \cdot (\nabla L)$	Computational efficient efficient	fixed ϵ .
SGD w/ Momentum	$\Delta w^i = \alpha \cdot \Delta w^{i-1} + (1-\alpha) \cdot (\nabla w^i)$	accelerate training	need hyperparameter α .
AdaGrad	$\Delta w^i = \frac{-\epsilon}{\sqrt{s + \sum (\nabla w^i)^2}} (\nabla w^i)$ $r_i = r^{i-1} + (\nabla L)^2$	No hyperp. Rescale ϵ adapt.	excessively decrease decrease ϵ .
RMS prop	$\Delta w^i = \frac{-\epsilon}{\sqrt{s + r_i}} (\nabla L)$ $r_i = \rho \cdot r^{i-1} + (1-\rho)(\nabla L)^2$	Discard extremely small gradients ↓ faster.	one additional ρ .
RMSprop with Nesterov M.	$\Delta w^i = \alpha \cdot \Delta w^{i-1} + \frac{\epsilon}{\sqrt{r_i}} (-\nabla L)$		α, ρ .
Adam.	$\Delta w^i = \frac{-\epsilon}{\sqrt{s + r_i}} \cdot \hat{s}^i$ $\hat{s}^i = \rho_1 \cdot s^{i-1} + (1-\rho_1)(\nabla L)$ $r_i = \rho_2 \cdot r^{i-1} + (1-\rho_2)(\nabla L)^2$	momentum bias correction	ρ_1, ρ_2