

Optimization Functions

27 May 2018 16:51

Parameter Update

Vanilla Update

Change the parameters along the negative gradient direction

$$w_i = w_i - \eta \frac{\partial C}{\partial w_i} \quad \text{where } \eta \text{ is learning Rate.}$$

A small η guarantees the progress in non-negative progress on loss function, towards minima. But small learning rate increases the training time.

Momentum Update

$$v = u \cdot \eta - \eta \frac{\partial C}{\partial w_i} \quad \text{where } u \text{ is momentum}$$

$$w_i = w_i + v$$

A typical momentum annealing schedule is to start with momentum of about 0.5 and anneal it to 0.99 or so over multiple epochs at later stages.

Nesterov Momentum

Stronger theoretical convergence guarantees for convex functions

$$w_{\text{ahead}} = w + u \times v$$

$$v = u \times v - \eta \times \frac{\partial C}{\partial w_{\text{ahead}}}$$

$$w = w + v$$

Rewriting the above equations (update in terms of w_{ahead} instead of w)

$$v_{\text{prev}} = v$$

$$v = u \times v - \eta \times \frac{\partial C}{\partial w}$$

$$w = w - u \times v_{\text{prev}} + (1 + u) \times v$$

9er-paramecer adapcive learning race mechods

Adagrad

$$C = C + dw^2$$

$$w = m - f_i \times \frac{dw}{\sqrt{C + \zeta}} \quad \zeta \text{ to avoid division by 0}$$

RMSprop

$$C = Q \times C + (1 - Q) \times dw^2 \quad \text{where } Q \text{ is decay rate}$$

$$w = w - \frac{f_i \times dw}{\sqrt{C + \zeta}}$$

Hence, RMSprop still modulates the learning rate of each weight based on the magnitudes of its gradients, which has a beneficial equalizing effect, but unlike Adagrad the updates do not get monotonically smaller.

Adam

$$m = Q_1 \times m + (1 - Q_1) \times dw \quad \# \text{ smoothed version of gradient}$$

$$v = Q_2 \times v + (1 - Q_2) \times dw^2$$

$$w = w - f_i \times \frac{m}{\sqrt{v + \zeta}} \quad \zeta \text{ to avoid division by 0}$$

Recommended Values of

$$\zeta = 1e - 8$$

$$Q_1 = 0.9$$

$$Q_2 = 0.999$$

Learning Rate Decay

- Step decay: Reduce the learning rate by some factor every few epochs. Typical values might be reducing the learning rate by a half every 5 epochs, or by 0.1 every 20 epochs.
- Exponential decay. has the mathematical form $a = a_0 e^{-kt}$, where a_0, k are hyper parameters and c is the iteration number (but you can also use units of epochs).
- 1/c decay has the mathematical form $a = \frac{a_0}{(1+kt)}$ where a_0, k are hyper parameters and c is the iteration number.

Source <<https://github.com/erogul/neural-networks-5/>>

