

Optimization Functions

27 May 2018 16:51

Parameter Update

Vanilla Update

Change the parameters along the negative gradient direction

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

A small λ guarantess the progress in non – negative progress on loss function, towards minima.
But small learning rate increases the traiging time.

Momentum Update

$$v = u \cdot v - \lambda \frac{\partial c}{\partial \omega_i} \quad \text{where } u \text{ is momentum}$$

$$w_i = w_i + v$$

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

Nestrov Momentum

Stronger theoretical converge guarantees for convex functions

$$w_{ahead} = w + u * v$$

$$v = u * v - \lambda * dw_{ahead}$$

$$w = w + v$$

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

$$v_{prev} = v$$

$$v = u * v - \lambda * dw$$

$$w = w - u * v_{prev} + (1 + u) * v$$

Per-parameter adaptive learning rate methods

Adagrad

$$C = C + dw^2$$

$$w = w - \lambda * \frac{dw}{\sqrt{C} + \epsilon} \quad \epsilon \text{ to avoid division by 0}$$

RMSprop

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

$$w = w - \frac{\lambda * dw}{\sqrt{C} + \epsilon}$$

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

$$m = \beta_1 * m + (1 - \beta_1) * dw \quad \# \text{ smooth version of gradient}$$

$$v = \beta_2 * v + (1 - \beta_2) * dw^2$$

$$w = w - \lambda * \frac{m}{\sqrt{v} + \epsilon} \quad \epsilon \text{ to avoid division by 0}$$

Recommended Values of

$$\epsilon = 1e-8$$

$$\beta_1 = 0.9$$

$$\beta_2 = 0.999$$

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 $\alpha = \alpha_0 e^{-kt}$, where α_0, k are hyper parameters and t is the iteration number (but you can also use units of epochs).
- **1/t decay** has the mathematical form $\alpha = \frac{\alpha_0}{(1+kt)}$ where α_0, k are hyper parameters and t is the iteration number.

