

Gradient descent algorithms

We use imperative assignment

$$\vec{\omega} := \vec{\omega} + \Delta \vec{\omega}$$

$\vec{\nabla} E(\vec{\omega})$ refers to the gradient of the error at $\vec{\omega}$

$\vec{\nabla}_i E(\vec{\omega})$ is the i -th component of the gradient, i.e. $\partial E / \partial \omega_i$

$\vec{\nabla}^{(i)} E(\vec{\omega})$ is the gradient computed for a single example (\vec{x}_i, \vec{y}_i)

BATCH

$$\vec{\omega} := \vec{\omega} - \alpha \vec{\nabla} E(\vec{\omega}) \quad \left. \begin{array}{c} \uparrow \\ \text{learning rate} \end{array} \right\} \text{one update every epoch}$$

On-line

$$\vec{\omega} := \vec{\omega} - \alpha \nabla^{(i)} E(\vec{\omega}) \quad \left. \right\} \text{one update for every example}$$

Mini-batch

$$\vec{\omega} := \vec{\omega} - \alpha \nabla^{(i:i+m)} E(\vec{\omega}) \quad \left. \right\} \text{one update per mini-batch}$$

Momentum

$$\Delta \vec{\omega} := \gamma \Delta \vec{\omega} - \eta \vec{\nabla} E(\vec{\omega})$$

$$\vec{\omega} := \vec{\omega} + \Delta \vec{\omega}$$

OR

$$\Delta \vec{\omega} := \beta \Delta \vec{\omega} + (1-\beta)(-\vec{\nabla} E(\vec{\omega}))$$

$$\vec{\omega} := \vec{\omega} + \Delta \vec{\omega}$$

If $\beta = 0.9$

$$\Delta \vec{\omega} := 0.9 \Delta \vec{\omega} + 0.1 (-\vec{\nabla} E(\vec{\omega}))$$

keep 90%
of last correction

add 10% of
the negative
gradient

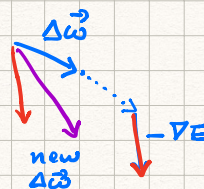
Momentum is only used with batch BP or minibatches, not with online BP

Nesterov accelerated gradient

$$\Delta \vec{\omega} := \gamma \Delta \vec{\omega} + \eta (-\vec{\nabla} E(\vec{\omega} + \gamma \Delta \vec{\omega}))$$

$$\vec{\omega} := \vec{\omega} + \Delta \vec{\omega}$$

If $\gamma = 0.9$, $\eta = 0.1$, keep 90% of last correction and add 10% of the negative gradient at an approximation of the gradient at the next point.



Adagrad

Adapt each weight with a different learning rate

sum of
squares of
past gradients →

$$G_i^t = G_i^{t-1} + (\Delta \omega_i)^2$$

$$\Delta \omega_i = - \frac{\alpha}{\sqrt{G_i^t + \epsilon}} \vec{\nabla}_i E(\vec{\omega})$$

the learning rate is smaller if the RMS of the past corrections is too big.

RMSprop

corrections are done element-wise

$$E[\Delta w_i^2] := \overset{\gamma}{0.9} E[\Delta w_i^2] + \overset{(1-\gamma)}{0.1} [\nabla_i E(\vec{w})^2]$$

$$\vec{w}_i := \vec{w}_i - \frac{\alpha}{\sqrt{E[\Delta w_i^2] + \epsilon}} \vec{\nabla}_i E(\vec{w})$$

Adam

$$\vec{m}_i := \beta_1 \vec{m}_i + (1-\beta_1) \vec{\nabla}_i E(\vec{w})$$

$$v_i := \beta_2 v_i + (1-\beta_2) \nabla_i E(\vec{w})^2$$

$$\vec{w} := \vec{w} - \frac{\alpha}{\sqrt{v} + \epsilon} \vec{m}$$

Two hyperparameters, one for tracking the mean, one for tracking the variance

(similar to the Kalman filter)