

Algorithm: SGD update at training iteration k

Require: Learning rate ϵ .

Require: Initial parameter θ

while stopping criterion not met do

Sample a minibatch of m examples from the training set $\{x^{(1)}, \dots, x^{(m)}\}$ with corresponding targets $y^{(i)}$.

Compute gradient estimate: $\hat{g} \leftarrow +\frac{1}{m} \nabla_{\theta} \sum_i L(f(x^{(i)}; \theta), y^{(i)})$

Apply update: $\theta \leftarrow \theta - \epsilon \hat{g}$

end while

$$\theta \leftarrow \theta - \epsilon (\hat{g})$$

1 Million
m

A weight with larger derivative will have larger updates

$$\epsilon_i \rightarrow w_i$$

(Actually we see the opposite)

$$\theta_i^{(k+1)} \leftarrow \theta_i^{(k)} - \epsilon (\nabla_i (\sum L(f(x_i^{(j)}; \theta), y)))$$

The learning parameter should not be constant for all iterations
(reason is out of scope)
also, from parameter to parameter

Adagrad

(Let's do this before Adam)

$\epsilon \rightarrow$ global
Initial θ

$$\delta = 10^{-7}$$

$r = 0$ (accumulation of gradient)

while

minibatch m

$$g \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(x_i), \theta, y^{(i)})$$

$$r \leftarrow r + g \odot g \quad \rightarrow \text{element wise product}$$

$$\Delta \theta \leftarrow -\frac{\epsilon}{\delta + \sqrt{r}} \odot g \quad \text{but } r \text{ is a vector}$$

$$\theta \leftarrow \theta + \Delta \theta$$

$$\frac{\epsilon}{\delta + \sqrt{r}}$$

$$\frac{\epsilon}{\delta + \sqrt{r}} = \begin{bmatrix} \frac{\epsilon}{\delta + \sqrt{r_1}} \\ \frac{\epsilon}{\delta + \sqrt{r_2}} \\ \frac{\epsilon}{\delta + \sqrt{r_3}} \\ \vdots \end{bmatrix}$$

$$\text{where } r = \begin{bmatrix} r_1 \\ r_2 \\ r_3 \\ \vdots \\ r_n \end{bmatrix}$$

L_1, L_2
Dropout
Early stopping
Data augmentation
multitasking
Adversarial training

RMSProp

- Modifies AdaGrad for a nonconvex setting
 - Change gradient accumulation into exponentially weighted moving average
 - Converges rapidly when applied to convex function

Geoffrey Hinton

Hinton

The RMSProp Algorithm

Require: Global learning rate ϵ , decay rate ρ .
 Require: Initial parameter θ
 Require: Small constant δ usually 10^{-9} , used to stabilize division by small numbers.
 Initialize accumulation variables $r = 0$
 while stopping criterion not met do
 Sample a minibatch of m examples from the training set $\{x^{(1)}, \dots, x^{(m)}\}$ with corresponding targets $y^{(i)}$.
 Compute gradient: $g \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(x^{(i)}; \theta), y^{(i)})$
 Accumulate squared gradient: $r \leftarrow \rho r + (1 - \rho) g \odot g$
 Compute parameter update: $\Delta \theta = -\frac{\epsilon}{\sqrt{r} + \delta} g$ (applied element-wise)
 Apply update: $\theta \leftarrow \theta + \Delta \theta$
 end while

0.6

$$r \leftarrow \rho r + (1 - \rho) (g \odot g)$$

BITS Pilani, Hyderabad Can

Algorithm: SGD with momentum

Require: Learning rate ϵ , momentum parameter α .
 Require: Initial parameter θ , initial velocity v .
 while stopping criterion not met do
 Sample a minibatch of m examples from the training set $\{x^{(1)}, \dots, x^{(m)}\}$ with corresponding targets $y^{(i)}$.
 Compute gradient estimate: $g \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(x^{(i)}; \theta), y^{(i)})$
 Compute velocity update: $v \leftarrow \alpha v - \epsilon g$
 Apply update: $\theta \leftarrow \theta + v$
 end while

→ We aren't comparing adagrad with momentum

→ SGD + momentum may perform better than adagrad at times

But it is better than SGD!

Algorithm: SGD with Nesterov momentum

Require: Learning rate ϵ , momentum parameter α .

Require: Initial parameter θ , initial velocity v .

while stopping criterion not met **do**

Sample a minibatch of m examples from the training set $\{x^{(1)}, \dots, x^{(m)}\}$ with corresponding labels $y^{(i)}$

Apply interim update: $\tilde{\theta} \leftarrow \theta + \alpha v$

This line is added from plain momentum

Compute gradient (at interim point): $g \leftarrow \frac{1}{m} \nabla_{\tilde{\theta}} \sum_i L(f(x^{(i)}; \tilde{\theta}), y^{(i)})$

Compute velocity update: $v \leftarrow \alpha v - \epsilon g$

Apply update: $\theta \leftarrow \theta + v$

end while

Algorithm: RMSProp with Nesterov momentum

Require: Global learning rate ϵ , decay rate ρ , momentum coefficient α .

Require: Initial parameter θ , initial velocity v .

Initialize accumulation variable $r = 0$

while stopping criterion not met **do**

Sample a minibatch of m examples from the training set $\{x^{(1)}, \dots, x^{(m)}\}$ with corresponding targets $y^{(i)}$.

Compute interim update: $\tilde{\theta} \leftarrow \theta + \alpha v$

Compute gradient: $g \leftarrow \frac{1}{m} \nabla_{\tilde{\theta}} \sum_i L(f(x^{(i)}; \tilde{\theta}), y^{(i)})$

Accumulate gradient: $r \leftarrow \rho r + (1 - \rho) g \odot g$

Compute velocity update: $v \leftarrow \alpha v - \frac{\epsilon}{\sqrt{r}} \odot g$. ($\frac{1}{\sqrt{r}}$ applied element-wise)

Apply update: $\theta \leftarrow \theta + v$

end while

Das has left the meeting

$$v \leftarrow \alpha v - \left[\frac{\epsilon}{\sqrt{r}} \odot g \right]$$

The Adam Algorithm

Require: Step size ϵ (Suggested default: 0.001)

Require: Exponential decay rates for moment estimates, ρ_1 and ρ_2 in $[0, 1)$.
(Suggested defaults: 0.9 and 0.999 respectively)

Require: Small constant δ used for numerical stabilization. (Suggested default: 10^{-8})

Require: Initial parameters θ

Initialize 1st and 2nd moment variables $s = 0, r = 0$

Initialize time step $t = 0$

while stopping criterion not met **do**

Sample a minibatch of m examples from the training set $\{x^{(1)}, \dots, x^{(m)}\}$ with corresponding targets $y^{(i)}$.

Compute gradient: $g \leftarrow \frac{1}{m} \nabla_{\theta} \sum_i L(f(x^{(i)}; \theta), y^{(i)})$

$t \leftarrow t + 1$

Update biased first moment estimate: $s \leftarrow \rho_1 s + (1 - \rho_1)g$

Update biased second moment estimate: $r \leftarrow \rho_2 r + (1 - \rho_2)g \odot g$

Correct bias in first moment: $\hat{s} \leftarrow \frac{s}{1 - \rho_1^t}$

Correct bias in second moment: $\hat{r} \leftarrow \frac{r}{1 - \rho_2^t}$

Compute update: $\Delta\theta = -\epsilon \frac{\hat{s}}{\sqrt{\hat{r} + \delta}}$ (operations applied element-wise)

Apply update: $\theta \leftarrow \theta + \Delta\theta$

end while

Influential
but sparse
features
give lower
updates

GRIP (CS7015)
to bring the moments as expected
(They are lower)
Since $s=r=0$

specific way of updating