# CME 216, ME 343 - Spring 2020 Eric Darve, ICME



Adagrad can yield nice acceleration in some cases but its formula to scale the learning rate has some issues:

$$s_i = \sum_k \left(rac{\partial L_k}{\partial x_i}
ight)^2$$

$$\Delta x_i = -rac{lpha}{\sqrt{s_i+\epsilon}} rac{\partial L_k}{\partial x_i}$$

If the convergence is slow, the gradient decays slowly.

As a result,  $s_i$  grows, which results in small steps and further slow down in convergence.

# nd further

To address this issue, RMSProp uses the idea of momentum so that the scaling factor has a more controlled growth:

$$s_i \leftarrow eta s_i + (1 - eta) \Big( rac{\partial L_k}{\partial x_i} \Big)^2$$

$$\Delta x_i = -rac{lpha}{\sqrt{s_i}+\epsilon} rac{\partial L_k}{\partial x_i}$$

# When $\beta$ is small, only recent values of the gradient will contribute.

When eta is close to 1,  $s_i$  is close to

$$pprox rac{1}{k_0} \sum_{s=0}^{k_0} \Big(rac{\partial L_{k-s}}{\partial x_i}\Big)^2$$

Another important element of RMSProp is that it helps around saddle points.

Near saddle points, the gradient  $\nabla_x L$  becomes very small. By scaling with  $1/(\sqrt{s_i} + \epsilon)$ , we can get a nice boost.

```
def RMSProp(W, s, lr, beta, batch_size):
    eps_stable = 1e-7
    g = W.grad / batch_size
    s = beta * s + (1-beta) * square(g) # element-wise square
    W -= lr * g / (sqrt(s) + eps_stable) # element-wise division
```



The last method we will cover is Adam. It addresses two limitations of RMSProp.

The first one is that it adds momentum to the gradient calculation as well.

So momentum is applied to two places: to advance the gradient squared and the gradient itself.

$$m_i \leftarrow eta_1 m_i + (1 - eta_1) rac{\partial L_k}{\partial x_i}$$

$$s_i \leftarrow eta_2 s_i + (1-eta_2) \Big(rac{\partial L_k}{\partial x_i}\Big)^2$$

The second modification has to do with initialization.

The problem with momentum is that before batch 0, all variables are initialized to 0.

0, all

As a result, we tend to get small values initially.

Basically initially we only see

$$m_i pprox (1-eta_1) rac{\partial L_k}{\partial x_i}$$

instead of

$$m_i pprox rac{\partial L_k}{\partial x_i}$$

This can be made more precise with the following derivation.

We will find how to rescale  $m_i$  to avoid this problem in the simple case where we assume that the gradient is constant.

It's not perfect but it will improve the initialization bias quite a bit.

Let's assume that the gradient is constant. In that case we just want

$$m_i = rac{\partial L_k}{\partial x_i}$$

at all steps.

14/28

Let's compare with what we are actually getting. Denote  $G_i$  the gradient and assume it is constant. $m_i \leftarrow eta_1 m_i + (1-eta_1)G_i$ 

If we solve this equation we get, at step k:

$$m_i^{(k)} = eta_1^k m_i^{(0)} + (1-eta_1) G_i \sum_{l=0}^{k-1} eta_1^l$$

# We chose $m_i^{(0)}=0.$ So, we get:



$$m_i^{(k)} = (1-eta_1^k)G_i$$

17/28

# The scaling factor we pick is therefore: $\dfrac{1}{1-eta_1^k}$

when processing batch k.

The final formulas are.

Gradient and gradient-squared:

$$m_i \leftarrow eta_1 m_i + (1 - eta_1) rac{\partial L_k}{\partial x_i}$$

$$s_i \leftarrow eta_2 s_i + (1 - eta_2) \Big( rac{\partial L_k}{\partial x_i} \Big)^2$$

20/28

Initialization bias:



21/28

# Gradient step:



22/28

 $\beta_1$  is for the gradient momentum. It is typically chosen equal to 0.9.

 $\beta_2$  is for the gradient squared momentum. This one is chosen very close to 1, 0.999. This means that  $s_i$  changes only very slowly.

The number of steps required for changes in  $s_i$  to be visible is on the order of 1,000 steps.

```
def adam(W, m, s, lr, batch_size, k):
    beta1 = 0.9
    beta2 = 0.999
   eps_stable = 1e-8
   g = W.grad / batch_size
   m = beta1 * m + (1. - beta1) * g
    s = beta2 * s + (1. - beta2) * square(g)
   m_hat = m / (1. - beta1 ** k)
    s_hat = s / (1. - beta2 ** k)
   W -= lr * m_hat / (sqrt(s_hat) + eps_stable)
```

# You can check these animations on the <u>CS231n</u> class page.



- SGD
- Momentum
- NAG
- Adagrad
- Adadelta
- Rmsprop



0.5

Momentum has a  $\beta$  close to 1, and tends to overshoot too much.

Adagrad and RMSProp scale the gradient near the saddle point which provides a nice boost.

# Additional reading

- <u>Adaptive Subgradient Methods for Online Learning and</u> Stochastic Optimization by Duchi et al., 2011
- <u>Adaptive Subgradient Methods for Online Learning and</u> Stochastic Optimization, International Symposium on Mathematical Programming 2012, by Duchi et al.

- <u>Qualitatively characterizing neural network optimization</u> problems by Goodfellow et al., ICLR, 2015
- <u>Practical Recommendations for Gradient-Based Training of</u> <u>Deep Architectures</u> by Bengio, 2013
- Efficient Backprop by LeCun et al., 2012
- ADADELTA: an adaptive learning rate method by Zeiler, 2012