

## 11.5. Accelerated Gradient Methods

Lecture based on “Dive into Deep Learning” <http://D2L.AI> (Zhang et al., 2020)

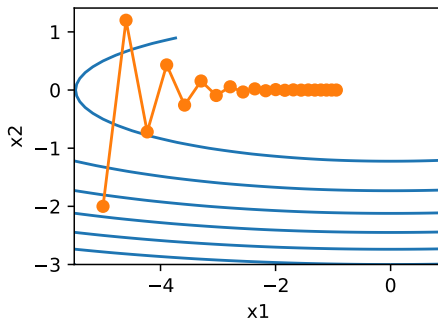
Prof. Dr. Christoph Lippert

Digital Health & Machine Learning

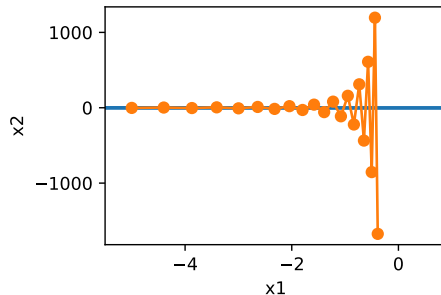
# Getting the learning rate right is crucial

## Example

$$f(x) = 0.1x_1^2 + 2x_2^2$$



learning rate 0.4.



learning rate 0.6.

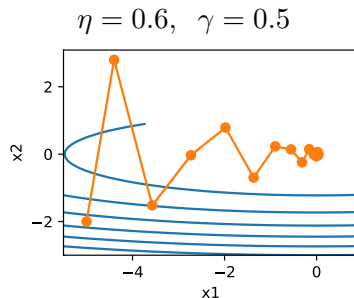
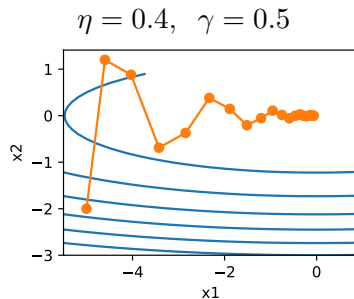
The **momentum** method solves this problem introducing a velocity  $\mathbf{v}_t$  at each step  $t$ .

$$\mathbf{v}_t \leftarrow \gamma \mathbf{v}_{t-1} + \eta_t \mathbf{g}_t,$$

$$\mathbf{x}_t \leftarrow \mathbf{x}_{t-1} - \mathbf{v}_t,$$

- At  $t = 0$ , momentum initializes the velocity variable  $\mathbf{v}_0 = \vec{0}$
- the **momentum hyperparameter**  $\gamma$  satisfies  $0 \leq \gamma < 1$ .

When  $\gamma = 0$ , momentum is equivalent to a mini-batch SGD.

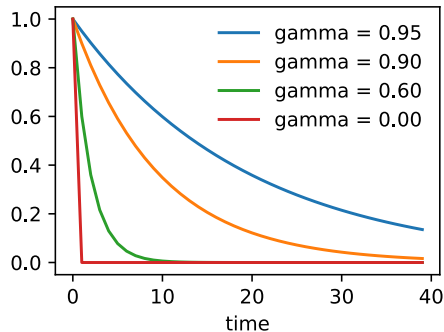


To understand momentum, we expand the velocity variable over time:

$$\begin{aligned}\mathbf{v}_t &= \eta_t \mathbf{g}_t + \gamma \mathbf{v}_{t-1}, \\ &= \eta_t \mathbf{g}_t + \gamma \eta_{t-1} \mathbf{g}_{t-1} + \gamma \mathbf{v}_{t-1}, \\ &\quad \dots \\ &= \eta_t \mathbf{g}_t + \gamma \eta_{t-1} \mathbf{g}_{t-1} + \dots + \gamma^{t-1} \eta_1 \mathbf{g}_1.\end{aligned}$$

- $\mathbf{v}_t$  is a weighted sum over all past gradients multiplied by the according learning rate
- which is the weight update in normal gradient descent.
- The weights decreases exponentially with the speed controlled by  $\gamma$ .

weights for 40 time steps under various  $\gamma$ s.



- Momentum smooths the weight updates over time.
- A small  $\gamma$  will focus on more recent scaled gradients.
- A large  $\gamma$  will include more past scaled gradients.
- It might walk out of a region that is flat.

Note that  $\frac{1}{1-\gamma} = 1 + \gamma + \gamma^2 + \dots$ .

- If all scaled gradients are similar to each other e.g.  $\eta_t \mathbf{g}_t \approx \eta \mathbf{g}$  for all  $t$ s,
- then the momentum changes the weight updates from  $\eta \mathbf{g}$  in gradient descent to  $\frac{\eta}{1-\gamma} \mathbf{g}$ .

## Momentum Summary

- The momentum method uses an exponentially weighted moving average (EWMA) of the gradient. It takes the weighted average of past time steps, with weights that decay exponentially by the time step.
- Momentum smooths independent variable updates for adjacent time steps in direction.

- In DL we apply a learning rate decay. (e.g.  $\eta_t = \frac{\eta_0}{\sqrt{t+c}}$ )
- If we observe gradient information in some dimensions only infrequently
  - e.g., in a language model where we observe some words only rarely (sparse features)
  - More generally, we have small partial derivatives in some dimensions most of the time and only rarely large ones
  - then the learning rate for this dimension might have already decayed too much for this dimension to still make progress
- Adagrad addresses this issue by tracking the aggregate of  $(\partial_i f(\mathbf{x}))^2$  for each dimension.

### Example

$f(\vec{x})$  with  $\vec{x} = [x_1, x_2]^\top$ , each element  $x_1$  and  $x_2$  uses the same learning rate  $\eta$ :

$$x_1 \leftarrow x_1 - \eta \frac{\partial f}{\partial x_1}, \quad x_2 \leftarrow x_2 - \eta \frac{\partial f}{\partial x_2}.$$

- When there is a big difference in gradient values for  $x_1$  and  $x_2$ , a small learning rate is needed to avoid divergence in larger gradient dimensions.
- This will cause the independent variables to move slowly in dimensions with smaller gradient values.

**Adagrad** adjusts the learning rate according to the gradient value in each dimension.



**Adagrad** re-adjusts the learning rate of each parameter using the cumulative variable  $\mathbf{s}_t$  capturing the square of the gradient dimensions:

$$\begin{aligned}\mathbf{s}_t &\leftarrow \mathbf{s}_{t-1} + \mathbf{g}_t \odot \mathbf{g}_t, \\ \mathbf{x}_t &\leftarrow \mathbf{x}_{t-1} - \frac{\eta}{\sqrt{\mathbf{s}_t + \epsilon}} \odot \mathbf{g}_t,\end{aligned}$$

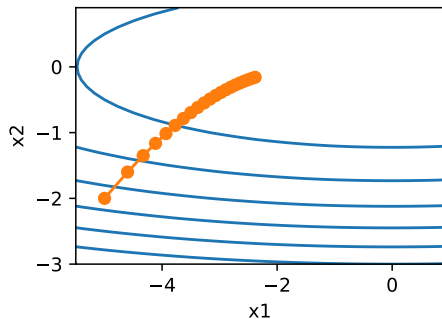
- $\odot$  is element-wise multiplication.
- The square root and division, and multiplication operations are all element operations.
- $\eta$  is the learning rate
- $\epsilon$  is a constant to maintain numerical stability, such as  $10^{-6}$ .
- At time step 0, Adagrad initializes each element in  $\mathbf{s}_0$  to 0.

$$\begin{aligned} \mathbf{s}_t &\leftarrow \mathbf{s}_{t-1} + \mathbf{g}_t \odot \mathbf{g}_t, \\ \mathbf{x}_t &\leftarrow \mathbf{x}_{t-1} - \frac{\eta}{\sqrt{\mathbf{s}_t + \epsilon}} \odot \mathbf{g}_t, \end{aligned}$$

- For elements that have a constant and large partial derivative, the learning rate will drop faster.
- If the partial derivative remains small, then its learning rate will decline more slowly.
- The learning rate of each element declines during each iteration.
- The learning rate declines very fast during an early iteration.
- Adagrad might have difficulty finding an optimum because the learning rate will be too small at later iterations.

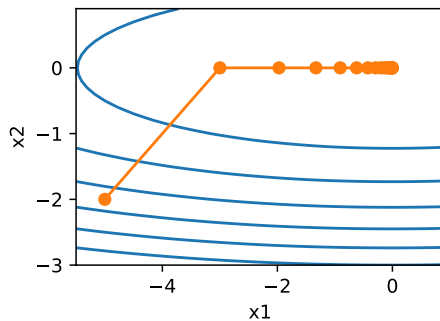
$$f(\mathbf{x}) = 0.1x_1^2 + 2x_2^2$$

learning rate 0.4



- the trajectory is smoother
- Due to the cumulative effect of  $s_t$ , the learning rate continuously decays

learning rate 2



- the independent variable approaches the optimum more quickly.

## Adagrad

# Summary

- Adagrad constantly adjusts the learning rate during each iteration to give each element in the independent variable of the objective function its own learning rate.
- When using Adagrad, the learning rate of each element in the independent variable decreases during each iteration.
- Adagrad might have difficulty finding a useful solution because the learning rate will be too small at later iterations.

In **RMSPprop** the entries of  $\mathbf{s}_t$  are computed as EWMA of the squares of the gradient elements.

$$\mathbf{s}_t \leftarrow \gamma \mathbf{s}_{t-1} + (1 - \gamma) \mathbf{g}_t \odot \mathbf{g}_t.$$

RMSPprop re-adjusts the learning rate as in Adagrad.

$$\mathbf{x}_t \leftarrow \mathbf{x}_{t-1} - \frac{\eta}{\sqrt{\mathbf{s}_t + \epsilon}} \odot \mathbf{g}_t,$$

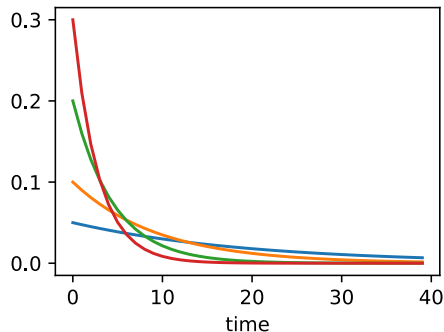
- $0 \leq \gamma < 1$  is a hyperparameter of the EWMA
- $\eta$  is the learning rate
- $\epsilon$  is a constant added to maintain numerical stability, such as  $10^{-6}$ .

Expanding the definition of  $\mathbf{s}_t$

$$\begin{aligned}\mathbf{s}_t &= (1 - \gamma)\mathbf{g}_t \odot \mathbf{g}_t + \gamma\mathbf{s}_{t-1} \\ &= (1 - \gamma)(\mathbf{g}_t \odot \mathbf{g}_t + \gamma\mathbf{g}_{t-1} \odot \mathbf{g}_{t-1}) + \gamma^2\mathbf{s}_{t-2} \\ &\dots \\ &= (1 - \gamma)(\mathbf{g}_t \odot \mathbf{g}_t + \gamma\mathbf{g}_{t-1} \odot \mathbf{g}_{t-1} + \dots + \gamma^{t-1}\mathbf{g}_1 \odot \mathbf{g}_1).\end{aligned}$$

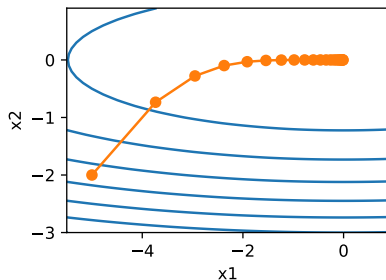
- $\frac{1}{1-\gamma} = 1 + \gamma + \gamma^2 + \dots$ , so the sum of weights equals to 1.
- The weights decrease exponentially

Weights in the past 40 time steps with  $\gamma \in [0.95, 0.9, 0.8, 0.7]$ .

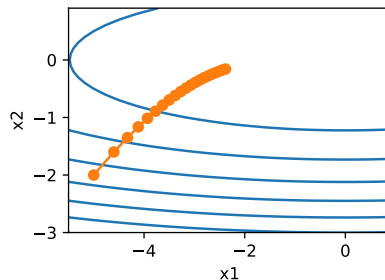


$$f(\mathbf{x}) = 0.1x_1^2 + 2x_2^2$$

$$\eta = 0.4, \gamma = 0.9$$



$$\text{Adagrad with } \eta = 0.4$$



- In contrast to Adagrad RMSProp can approach the optimal solution faster at  $\eta = 0.4$ .
- The difference between **RMSProp** and **Adagrad** is that RMSProp uses an EWMA on the squares of elements to adjust the learning rate.

**Adadelta** keeps track of two EWMA's,  $s_t$  and  $\Delta x_t$ .

$$s_t \leftarrow \rho s_{t-1} + (1 - \rho) g_t \odot g_t.$$

$\Delta x_{t-1}$  is used to compute the step  $g'_t$ :

$$g'_t \leftarrow \sqrt{\frac{\Delta x_{t-1} + \epsilon}{s_t + \epsilon}} \odot g_t,$$

$$x_t \leftarrow x_{t-1} - g'_t.$$

$\Delta x$  records the EWMA on the squares of the steps  $g'$

$$\Delta x_t \leftarrow \rho \Delta x_{t-1} + (1 - \rho) g'_t \odot g'_t.$$

- $0 \leq \rho < 1$  (counterpart of  $\gamma$  in RMSProp)
- $\epsilon$  is a constant added to maintain the numerical stability, such as  $10^{-5}$ .
- Adadelta replaces the learning rate  $\eta$  in RMSProp with  $\sqrt{\Delta x_{t-1}}$ .



**Adam** uses the **momentum**  $\mathbf{v}_t$  as the EWMA over the gradient.

$$\mathbf{v}_t \leftarrow \beta_1 \mathbf{v}_{t-1} + (1 - \beta_1) \mathbf{g}_t.$$

$\mathbf{s}_t$  is the EWMA on the squares of elements in the gradient

$$\mathbf{s}_t \leftarrow \beta_2 \mathbf{s}_{t-1} + (1 - \beta_2) \mathbf{g}_t \odot \mathbf{g}_t.$$

- $0 \leq \beta_1 < 1$  (the author suggests 0.9),
- $0 \leq \beta_2 < 1$  (the author suggests 0.999)

Because of initialization,  $\mathbf{v}_0 = 0$  and  $\mathbf{s}_0 = 0$ ,

$$\mathbf{v}_t = (1 - \beta_1) \sum_{i=1}^t \beta_1^{t-i} \mathbf{g}_i$$

For example, when  $\beta_1 = 0.9$ ,  $\mathbf{v}_1 = 0.1\mathbf{g}_1$ .

- The sum of weights from all previous time steps is  $(1 - \beta_1) \sum_{i=1}^t \beta_1^{t-i} = 1 - \beta_1^t$ .
- For small  $t$ , the sum will be small.

Adam corrects for this bias in  $\mathbf{v}_t$  and  $\mathbf{s}_t$ :

$$\hat{\mathbf{v}}_t \leftarrow \frac{\mathbf{v}_t}{1 - \beta_1^t}, \quad \hat{\mathbf{s}}_t \leftarrow \frac{\mathbf{s}_t}{1 - \beta_2^t}.$$

Adam uses the bias-corrected variables  $\hat{\mathbf{v}}_t$  and  $\hat{\mathbf{s}}_t$  to re-adjust the learning rate of each model parameter

$$\mathbf{g}'_t \leftarrow \frac{\eta \hat{\mathbf{v}}_t}{\sqrt{\hat{\mathbf{s}}_t} + \epsilon},$$

- $\eta$  is the learning rate
- $\epsilon$  is a constant added to maintain numerical stability, such as  $10^{-8}$ .

Finally, use  $\mathbf{g}'_t$  to update the parameters:

$$\mathbf{x}_t \leftarrow \mathbf{x}_{t-1} - \mathbf{g}'_t.$$

## Summary

- Created on the basis of RMSProp, Adam also uses EWMA on the mini-batch stochastic gradient
- Adam uses bias correction.

## Summary

- Most obviously the **magnitude** of the learning rate matters.
  - If it is too large, optimization diverges
  - If it is too small, it takes too long to train or we end up with a suboptimal result.
- The rate of decay is just as important.
  - If the learning rate remains large we may simply end up bouncing around the minimum and thus not reach optimality.
  - In short, we want the rate to decay, but probably more slowly than  $\mathcal{O}(t^{-\frac{1}{2}})$  which would be a good choice for convex problems.
- Another aspect that is equally important is **initialization**. This pertains both to how the parameters are set initially and also how they evolve initially.