

# EMAT31530: Practical optimization in AI

Laurence Aitchison

Now, we know how PyTorch computes gradients. However, to practically optimize neural networks, you need a bit more than just gradients: and that’s what we’ll look at this week!

## 1 Minibatching

So far, we’ve computed the full gradient based on all the training data in one big batch computation. In deep learning, this is typically called a “full batch” method. However, remember that GPU memory is very limited, and training datasets are often very, very large. So full batch training only works in small-scale settings.

As an alternative, we usually use minibatches. A minibatch is a small subset of the training data (e.g. 128 images out of 60,000 images in the full dataset). In minibatched training, we calculate the loss, compute the gradients and update the parameters for only the minibatch.

To write down minibatched updates, we need notation for the gradient for the  $i$ th datapoint,

$$\mathbf{g}_i = \frac{d\mathcal{L}_i}{d\mathbf{w}}. \quad (1)$$

Here,  $\mathcal{L}_i$  is the loss for the  $i$ th datapoint, and  $\mathbf{g}_i$  is the gradient for the  $i$ th datapoint. We can use the individual datapoint gradient to write down the minibatched gradient. The minibatched gradient is just the gradient, averaged over datapoints in the minibatch,

$$\mathbf{g}_{\text{mb}} = \frac{1}{M} \sum_{i \text{ in mb}} \mathbf{g}_i = \frac{1}{M} \sum_{i \text{ in mb}} \frac{d\mathcal{L}_i}{d\mathbf{w}}. \quad (2)$$

Here,  $M$  is the minibatch size, and “ $i$  in mb” gives the set of indices of datapoints in the current minibatch. However, it is quite awkward to use this expression directly in PyTorch, as, taken literally, it seems to involve backpropagating through  $M$  single-datapoint losses,  $\mathcal{L}_i$ . In contrast, PyTorch only allows us to (easily) work with a single loss. We therefore need to define a minibatched

loss. Note that we again use an *average*, rather than a sum, as it makes the code + maths slightly easier.

$$\mathcal{L}_{\text{mb}} = \frac{1}{M} \sum_{i \text{ in mb}} \mathcal{L}_i. \quad (3)$$

Using PyTorch to differentiate the minibatch loss gives the minibatched gradient,  $\mathbf{g}_{\text{mb}}$ ,

$$\frac{d\mathcal{L}_{\text{mb}}}{d\mathbf{w}} = \frac{d}{d\mathbf{w}} \left( \frac{1}{M} \sum_{i \text{ in mb}} \mathcal{L}_i \right) = \frac{1}{M} \sum_{i \text{ in mb}} \frac{d\mathcal{L}_i}{d\mathbf{w}} = \frac{1}{M} \sum_{i \text{ in mb}} \mathbf{g}_i = \mathbf{g}_{\text{mb}}. \quad (4)$$

The simplest optimizer that uses minibatching is called “stochastic gradient descent”, commonly abbreviated to SGD. This is very similar to gradient descent that we’ve seen previously, but it uses minibatched gradients, rather than the full batch gradient,

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \mathbf{g}_{\text{mb};t}, \quad (5)$$

where  $\mathbf{g}_{\text{mb};t}$  is the minibatched gradient at time  $t$ . Note that the “stochastic” in SGD refers to the stochasticity induced by randomness in the choice of minibatches.

This additional stochasticity raises a question: what theory do we have to say that SGD still “works”? And what do we even mean by “works” in this context? Well, a pretty important feature of an algorithm is whether it eventually converges to the minimum of the loss. That will (assuming there are no instabilities) happen for e.g. SGD applied in the full-batch setting with a fixed learning rate. However, if you use a fixed learning rate, SGD will get close to the optimum, but won’t actually exactly converge to exactly the optimum. Instead, it bounce around the optimum forever due to the “noisy” minibatched gradient estimates.

However, it turns out we can prove convergence in a slightly modified setting: as the learning rate,  $\eta$ , decays towards. Specifically, we can show that if you reduce the learning rate  $\eta$ , and at the same time, consider an increasing number of iterations,  $T(\eta) = 1/\eta$ , then you expect the SGD updates to move the parameters to the same place, but with decreasing variance (the variance is induced by stochasticity in the choice of minibatches). Reducing the learning rate as we increase the number of timesteps in effect implies we are averaging over more samples of the minibatch noise.

We start by proving that the expected minibatched update is the same as the full-batch update. Specifically, we’re going to write the full batch loss again as an average,

$$\mathcal{L}_{\text{fb}} = \frac{1}{N} \sum_{i=1}^N \mathcal{L}_i. \quad (6)$$

where  $N$  is the number of datapoints in the full dataset. The full batch gradient,  $\mathbf{g}_{\text{fb}}$  is again the average of individual datapoint gradients,

$$\mathbf{g}_{\text{fb}} = \frac{d\mathcal{L}_{\text{fb}}}{d\mathbf{w}} = \frac{1}{N} \sum_{i=1}^N \frac{d\mathcal{L}_i}{d\mathbf{w}} = \frac{1}{N} \sum_{i=1}^N \mathbf{g}_i. \quad (7)$$

Now, we can think of the individual datapoint gradient as the random variable you get by randomly selecting a datapoint, then computing its expectation. Critically, this random variable has the right expectation,

$$\mathbb{E}[\mathbf{g}_{\text{sd}}] = \frac{1}{N} \sum_{i=1}^N \mathbf{g}_i = \mathbf{g}_{\text{fb}}. \quad (8)$$

And as the minibatched gradient is just the average of a bunch of independent realisations of  $\mathbf{g}_{\text{sd}}$  (Eq. 2),

$$\mathbb{E}[\mathbf{g}_{\text{mb}}] = \mathbf{g}_{\text{fb}}, \quad (9)$$

it also has the right expectation.

Thus, we interpreted the single-datapoint,  $\mathbf{g}_{\text{sd}}$ , or minibatched,  $\mathbf{g}_{\text{mb}}$ , gradients as random variables, and proved that they have the right expectation,  $\mathbf{g}_{\text{fb}}$ . A statistician would say that  $\mathbf{g}_{\text{sd}}$  and  $\mathbf{g}_{\text{mb}}$  are “unbiased estimators” of  $\mathbf{g}_{\text{fb}}$ , because  $\mathbb{E}[\mathbf{g}_{\text{sd}}] = \mathbf{g}_{\text{fb}}$  and  $\mathbb{E}[\mathbf{g}_{\text{mb}}] = \mathbf{g}_{\text{fb}}$ . Now, to get a really good estimate of the true-batch gradient,  $\mathbf{g}_{\text{fb}}$ , all we need to do is to average over many independent realisations of  $\langle g \rangle$  or  $\mathbf{g}_{\text{mb}}$ . And it turns out we can do this implicitly by considering shrinking the learning rate as we increase the number of timesteps. Specifically, we consider a number of timesteps that is inversely proportional to the learning rate,

$$T(\eta) = \frac{1}{\eta}. \quad (10)$$

Note that we make sure to choose values of  $\eta$  that give integer numbers of steps  $T(\eta)$ . The overall change in the parameter,  $w$  is,

$$\Delta w = \eta \sum_{t=1}^{T(\eta)} g_{\text{mb};t}. \quad (11)$$

Writing  $\eta = 1/T(\eta)$ ,

$$\Delta w = \frac{1}{T(\eta)} \sum_{t=1}^{T(\eta)} g_{\text{mb};t} \quad (12)$$

Thus,  $\Delta w$  is just the average of  $T(\eta)$  independent realisations of the minibatched gradient,  $g_{\text{mb};t}$ . Assuming that we consider small changes in  $w$ , such that the expected gradient doesn’t change much across the small part of the trajectory

considered (i.e.  $E[g_{\text{mb};t}] = g_{\text{fb}}$ , where  $g_{\text{fb}}$  doesn't change over timesteps), the expected change in weights becomes,

$$E[\Delta w] = \frac{1}{T(\eta)} \sum_{t=1}^{T(\eta)} E[g_{\text{mb};t}] = \frac{1}{T(\eta)} \sum_{t=1}^{T(\eta)} g_{\text{fb}} = g_{\text{fb}}, \quad (13)$$

which is just the full-batch gradient. Moreover, as  $\Delta w$  is the average of  $T(\eta)$  independent realisations of the minibatch gradient,  $g_{\text{mb};t}$ , the variance of  $\Delta w$  is proportional to  $1/T(\eta)$ , or to  $\eta$  (Eq. 10),

$$\text{Var}[\Delta w] = \frac{1}{T(\eta)} \text{Var}[g_{\text{mb}}(t)] = \eta \text{Var}[g_{\text{mb}}(t)], \quad (14)$$

Thus, as  $\eta \rightarrow 0$ , the variance in the change in weights  $\text{Var}[\Delta w]$  also goes to zero.

Theory is great, but deep learning is a practical subject, so there's a few points of "wisdom" and/or common practice. Specifically:

- Each pass through the full training dataset is called an "epoch". Deep learning training progress is typically measured in epochs. Each epoch is composed of many minibatches (minibatches per epoch = dataset size / minibatch size).
- Within each epoch, it works slightly better to randomise exactly which datapoints are part of each minibatch, rather than going through the same fixed set of minibatches at each epoch again and again (that's what `shuffle=True` does).

## 2 The problem with NN optimization

The simple form of SGD described above doesn't work well, and is rarely used in practice. The problem all comes down to the interactions between the size of gradients and the learning rates. Specifically:

- If all the gradients are small, then we'd like to use large learning rates to converge reasonably quickly.
- If all the gradients are large, then we'd like to use small learning rates to avoid instability.

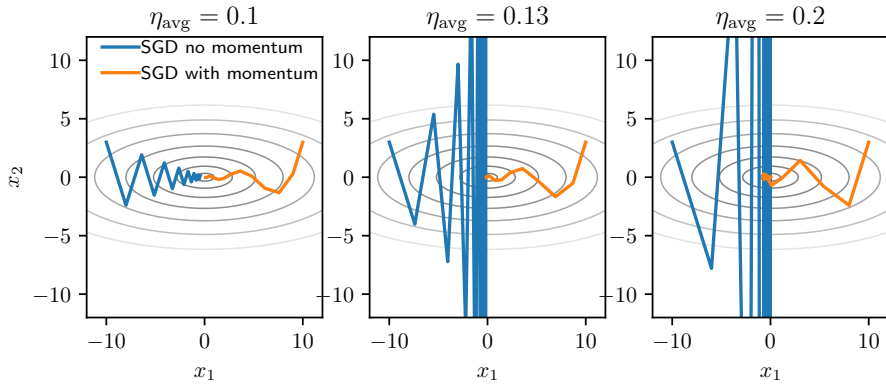
The problem is in practice, gradients in some directions are large, and gradients in some directions are small. That sets up a conflict: for directions with a large gradient, you'd like to use a small learning rate (to avoid instability). But for directions with a small gradient, you'd like to use a large learning rate (to converge reasonably quickly).

In practice, real problems have some small and some large gradients. In that case, you're forced to use a low learning rate (to avoid instability along directions with large gradients), which implies slow convergence along directions with small

gradients. It is this problem that all the optimization tricks discussed below are trying to address.

### 3 Momentum

Momentum is one way to mitigate these issues, as it stabilises learning with large-ish learning rates along high gradient directions. That allows you to push up the learning rate a bit, allowing faster convergence in the small-gradient directions. Specifically, the instabilities often look like the diagram below,



SGD with no momentum is plotted in the blue lines. If the learning rate is large, SGD with no momentum can oscillate: at each step, it jumps over the “horizontal” in these diagrams. That’s okay if the oscillations decay (i.e. left plot below with  $\eta = 0.1$ ), in which case you will still converge. But as the learning rate increases, the oscillations grow, and the learning process diverges (middle and right plots). Momentum (orange line) allows you to increase learning rates a bit while avoiding this instability.

How does momentum work? One intuition for momentum (and the origin of the term “momentum”) is the analogy with Physics. Specifically, consider a ball rolling around on a hilly landscape, with friction. The objective,  $\mathcal{L}(\mathbf{w})$  describes the height of the landscape at any location,  $\mathbf{w}$ . “Gravity” pushes the ball downhill with a force given by the gradient

$$\mathbf{F}(\mathbf{w}) = -\eta_{\text{trad}} \mathbf{g}_{\text{fb}}. \quad (15)$$

(Note that I’m using two slightly different versions of the momentum update, with two slightly different learning rates: a traditional formulation with learning rate,  $\eta_{\text{trad}}$ , that you’ll find on the internet, and an “average” formulation with learning rate  $\eta_{\text{avg}}$  that makes the effect of momentum and the connection to later methods such as Adam far clearer). Taking the mass and gravitational

acceleration to be 1 for simplicity, the equations of motion are,

$$\frac{d\mathbf{v}}{dt} = -(1 - \mu)\mathbf{v} + \mathbf{F}(\mathbf{w}) = -(1 - \mu)\mathbf{v} - \eta_{\text{trad}}\mathbf{g}_{\text{fb}} \quad (16a)$$

$$\frac{d\mathbf{w}}{dt} = \mathbf{v}. \quad (16b)$$

Here, the  $-(1 - \mu)\mathbf{v}$  term in the first equation represents friction ( $\mu$  is between 0 and 1). Discretising these equations with a timestep of 1 and using the mini-batched gradient,  $\mathbf{g}_{\text{mb};t}$  in place of the full batch gradient,  $\mathbf{g}_{\text{fb}}$ , gives

$$\mathbf{v}_{t+1} = \mathbf{v}_t + \frac{d\mathbf{v}}{dt} = \mu\mathbf{v}_t - \eta_{\text{trad}}\mathbf{g}_{\text{mb};t} \quad (17a)$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \frac{d\mathbf{w}}{dt} = \mathbf{w}_t + \mathbf{v}_{t+1}. \quad (17b)$$

This is the traditional parametrisation e.g. used in PyTorch (SGD docs). Typically we use a momentum of  $\mu = 0.9$ , while the learning rate,  $\eta_{\text{trad}}$ , needs extensive tuning.

However, I’m really not a fan of the traditional parameterisation: it doesn’t help us to understand what’s really going on, and conflicts with how other learning rules such as Adam are written down. These other learning rules use an exponential moving average gradient,

$$\langle \mathbf{g} \rangle_{t+1} = \beta_1 \langle \mathbf{g} \rangle_t + (1 - \beta_1)\mathbf{g}_{\text{mb};t} \quad (18)$$

This looks alot the “velocity” updates in Eq. (17a), in that they both have an exponentially decay, with an additive term that depends on the minibatch gradient (the exponential decay is for  $\mathbf{v}$  in Eq. (17a) and for  $\langle \mathbf{g} \rangle$  in Eq. (18)). However, this form is much easier to interpret, as it is an exponential moving average of minibatch gradients that just estimates the full batch gradient (see exercises for details). That’s why we denote this quantity  $\langle \mathbf{g} \rangle_t$  (angle brackets are often used to denote some type of average). Now, we can write an alternative set of momentum updates,

$$\langle \mathbf{g} \rangle_{t+1} = \beta_1 \langle \mathbf{g} \rangle_t + (1 - \beta_1)\mathbf{g}_{\text{mb};t} \quad (19a)$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_{\text{avg}} \langle \mathbf{g} \rangle_{t+1} \quad (19b)$$

This is alot more similar to traditional SGD, in that the parameter updates are again just a learning rate times an estimate of the gradient. In SGD, we just use a single minibatched gradient estimate, while here, we use an exponential-moving-average over previous minibatch gradient estimates.

This notion of the exponential-moving-averaging gradients allows us to understand why momentum stabilises learning by suppressing oscillations. Specifically, what happens if we do an exponential moving average over gradients that alternate sign? Well, the positive and negative terms in that exponential moving average are going to partially cancel, leading to an average gradient

that's much smaller in magnitude in these oscillating directions. And that much smaller gradient implies smaller updates and improved stability.

The two formulations are equivalent if we set:

$$\mu = \beta_1 \qquad \eta_{\text{avg}} = \frac{1}{1-\beta_1} \eta_{\text{trad}}. \quad (20)$$

That means if we increase momentum (typically we use  $\mu = \beta_1 = 0.9$ ) while keeping the traditional momentum,  $\eta_{\text{trad}}$ , the same, then we are in effect increasing the learning rate,  $\eta_{\text{avg}}$ .

### 3.1 Proof of relationship between the two formulations [non-examinable]

Now, the existence of these two parameterisations raises the question of how they interrelate. Specifically, we want to find relationships between  $\mu$  and  $\eta_{\text{trad}}$  in the traditional formulation and  $\beta_1$  and  $\eta_{\text{avg}}$  in the average formulation such that the traditional and average formulations give exactly the same parameter updates at all timesteps. Specifically, we need the parameter update from the traditional formulation (Eq. 17b) to equal to the update from the alternative formulation (Eq. 19b) at timestep  $t$  and  $t + 1$ ,

$$\mathbf{v}_t = -\eta_{\text{avg}} \langle \mathbf{g} \rangle_t \quad (21a)$$

$$\mathbf{v}_{t+1} = -\eta_{\text{avg}} \langle \mathbf{g} \rangle_{t+1}. \quad (21b)$$

This indicates that we need the momentum in the traditional formulation to be directly proportional to the average gradient. To achieve that, let's substitute the updates for  $\mathbf{v}_{t+1}$  (Eq. 17a) and for  $\langle \mathbf{g} \rangle_{t+1}$  (Eq. 19a) into Eq. (21b),

$$\begin{aligned} \mu \mathbf{v}_t - \eta_{\text{trad}} \mathbf{g}_{\text{mb};t} &= -\eta_{\text{avg}} (\beta_1 \langle \mathbf{g} \rangle_t + (1 - \beta_1) \mathbf{g}_{\text{mb};t}) \\ &= -\eta_{\text{avg}} \beta_1 \langle \mathbf{g} \rangle_t - \eta_{\text{avg}} (1 - \beta_1) \mathbf{g}_{\text{mb};t}. \end{aligned} \quad (22)$$

Now, let's use (Eq. 21a) to substitute for  $\mathbf{v}_t$  on the LHS,

$$-\mu \eta_{\text{avg}} \langle \mathbf{g} \rangle_t - \eta_{\text{trad}} \mathbf{g}_{\text{mb};t} = -\eta_{\text{avg}} \beta_1 \langle \mathbf{g} \rangle_t - \eta_{\text{avg}} (1 - \beta_1) \mathbf{g}_{\text{mb};t}. \quad (23)$$

Now, the minibatch gradient,  $\mathbf{g}_{\text{mb};t}$  is a random variable that could (in principle) be anything, and is not e.g. deterministically tied to  $\langle \mathbf{g} \rangle_t$ . Thus, we're going to consider the equality of the terms proportional to the minibatch gradient, and the terms proportional to  $\langle \mathbf{g} \rangle_t$  separately. For the terms proportional to  $\langle \mathbf{g} \rangle_t$  to be equal, we need,

$$-\mu \eta_{\text{avg}} \langle \mathbf{g} \rangle_t = -\eta_{\text{avg}} \beta_1 \langle \mathbf{g} \rangle_t. \quad (24)$$

Thus,

$$\mu = \beta_1. \quad (25)$$

And for the terms proportional to  $\mathbf{g}_{\text{mb};t}$  to be equal,

$$-\eta_{\text{trad}}\mathbf{g}_{\text{mb};t} = -\eta_{\text{avg}}(1 - \beta_1)\mathbf{g}_{\text{mb};t}. \quad (26)$$

Thus,

$$\eta_{\text{trad}} = \eta_{\text{avg}}(1 - \beta_1) \quad \eta_{\text{avg}} = \frac{1}{1 - \beta_1}\eta_{\text{trad}}. \quad (27)$$

Which gives our result.

## 4 RMSProp

If you remember, the original problem we were solving with momentum was that there are big gradients in some directions, and small gradients in other directions. You ideally want a big learning rate in directions with a small gradient, to converge quickly. At the same time, you want a small learning rate in directions with a big gradient, to avoid instability. Momentum mitigated the problem, by stabilising along high-gradient directions. But it is also possible to just directly reduce the gradients in the high-gradient directions. And that's just what RMSProp gives you! Specifically, RMSProp computes an exponential moving average of the squared gradient,  $g_{\text{mb};t}^2$  then uses the root-mean-square-average gradient to normalize (or adapt) the learning rate,

$$\langle g^2 \rangle_{t+1} = \beta_2 \langle g^2 \rangle_t + (1 - \beta_2)g_{\text{mb};t}^2 \quad (28)$$

$$w_{t+1} = w_t + \frac{\eta}{\sqrt{\langle g^2 \rangle_{t+1} + \epsilon}} g_{\text{mb};t}. \quad (29)$$

Note that  $\epsilon$  is a small positive constant (usually set to  $10^{-8}$ ) which ensures that we never divide-by-zero.

To understand how RMS prop functions, we assume that we're in a small region, where the gradient is always the same,  $g_{\text{mb};t} = g_{\text{fb}}$ , and if we're in this region for long enough,  $\langle g^2 \rangle_{t+1} = g_{\text{fb}}^2$ . Now, the weight update is,

$$w_{t+1} = w_t + \eta \frac{g_{\text{fb}}}{\sqrt{g_{\text{fb}}^2 + \epsilon}} \quad (30)$$

And  $\sqrt{g_{\text{fb}}^2} = |g_{\text{fb}}|$ ,

$$w_{t+1} = w_t + \eta \frac{g_{\text{fb}}}{|g_{\text{fb}}| + \epsilon} \quad (31)$$

Setting  $\epsilon = 0$  for simplicity,

$$w_{t+1} = w_t + \eta \text{sign}(g_{\text{fb}}) \quad (32)$$

$$\text{sign}(x) = \begin{cases} 1 & \text{if } 0 < x \\ 0 & \text{if } 0 = x \\ -1 & \text{if } 0 > x \end{cases} \quad (33)$$



So the magnitude of the updates is always around  $\eta$ , no matter how big the gradients are. This turns out to be super-important (but not widely understood) intuition for setting  $\eta$  to sensible values and for setting up your model to work well with Adam/RMSProp.

## 5 Adam

In practice, RMSProp is rarely, if ever used. Instead, Adam is perhaps the most common optimizer in deep learning. Critically, Adam is just RMSProp + Momentum + Bias correction. Note that the name “Adam” breaks down as “Ada-m”, with the “Ada” for “adaptive” learning rates, and the “m” for momentum. To combine RMSProp with momentum, we just replace the gradient,  $g_{\text{mb};t}$ , in Eq. (29) with the exponential moving average gradient,

$$\langle g \rangle_{t+1} = \beta_1 \langle g \rangle_t + (1 - \beta_1) g_{\text{mb};t} \quad (34)$$

$$\langle g^2 \rangle_{t+1} = \beta_2 \langle g^2 \rangle_t + (1 - \beta_2) g_{\text{mb};t}^2 \quad (35)$$

$$w_{t+1} = w_t + \frac{\eta}{\sqrt{\langle g^2 \rangle_{t+1} + \epsilon}} \langle g \rangle_{t+1}. \quad (36)$$

These equations show Adam updates without debiasing. However, there’s a problem with these updates: specifically,  $\langle g^2 \rangle_{t+1}$  is too small for about the first 1000 timesteps. In particular, we typically use a long timescale for the exponential moving average for  $\langle g^2 \rangle_{t+1}$ : the default is  $\beta_2 = 0.999$ , which implies a timescale of about 1000 timesteps. Moreover, in the absence of any better information, we initialize  $\langle g^2 \rangle_0 = 0$ . Combining all this, we have  $1 - \beta_2 = \frac{1}{1000}$ , so

$$\langle g^2 \rangle_1 = \frac{1}{1000} g_{\text{mb};t=1}^2, \quad (37)$$

thus the first estimate of  $\langle g^2 \rangle$  is about 1000 times too small. Over about 1000 timesteps,  $\langle g^2 \rangle$  will converge to more sensible values. But it still will have spent a long time with a value that is too high. To fix this issue, the full Adam updates do bias correction:

$$\langle g \rangle_{t+1} = \beta_1 \langle g \rangle_t + (1 - \beta_1) g_{\text{mb};t} \quad (38)$$

$$\langle g^2 \rangle_{t+1} = \beta_2 \langle g^2 \rangle_t + (1 - \beta_2) g_{\text{mb};t}^2 \quad (39)$$

$$\hat{m}_{t+1} = \frac{\langle g^2 \rangle_{t+1}}{1 - (\beta_1)^t} \quad (40)$$

$$\hat{v}_{t+1} = \frac{\langle g \rangle_{t+1}}{1 - (\beta_2)^t} \quad (41)$$

$$w_{t+1} = w_t + \frac{\eta}{\sqrt{\hat{v}_{t+1} + \epsilon}} \hat{m}_{t+1}. \quad (42)$$

To understand what bias correction is doing, see Exercises.

## 6 Gradient clipping

The key problem inducing instability was that the gradients in some directions were too large. Clipping is perhaps the simplest way to prevent this problem. Usually, the clip function is defined as,

$$\text{clip}(x, \max, \min) = \begin{cases} \max & \text{if } \max < x \\ \min & \text{if } x < \min \\ x & \text{otherwise} \end{cases} \quad (43)$$

In gradient clipping, we'd apply clip to the gradients, with  $\min = -\max$ .

Gradient clipping seems to work reasonably well in practice, and is often applied. However, there is an important health-warning when applying gradient clipping in the context of minibatched optimization. Remember that a key property that we need for optimization to work (e.g. to converge to the right answer) was that the expected minibatch gradient should be equal to the full batch gradient,

$$\mathbb{E}[\mathbf{g}_{\text{mb}}] = \mathbf{g}_{\text{fb}}. \quad (44)$$

The problem is that the expected *clipped* gradient is no longer equal to the full batch gradient,

$$\mathbb{E}[\text{clip}(\mathbf{g}_{\text{mb}}, \min, \max)] \neq \mathbf{g}_{\text{fb}}, \quad (45)$$

because clip is a nonlinear function. This might not be a problem if it happens rarely, or if it happens increasingly rarely if the model converges. But it is always worth bearing in mind when using clipping!

## 7 Exercises

**Exercise 1.** Consider an exponential moving average of random variables  $X_{t+1}$ ,

$$\langle X \rangle_{t+1} = \beta \langle X \rangle_t + (1 - \beta) X_{t+1} \quad (46)$$

where we initialize at,

$$\langle X \rangle_0 = 0 \quad (47)$$

and the means of all the  $X_{t+1}$ 's are  $\mu$ ,

$$\mathbb{E}[X_{t+1}] = \mu. \quad (48)$$

Part 1: What is  $\mathbb{E}[\langle X \rangle_{t+1}]$  after an arbitrary number of timesteps?

Part 2: Can I multiply by a time-dependent constant,  $c_{t+1}$  to form a “debiased” estimate,

$$d_{t+1} = c_{t+1} \langle X \rangle_{t+1} \quad (49)$$

such that

$$\mathbb{E}[d_{t+1}] = \mu. \quad (50)$$

**Answer 1.** *Part 1: We begin by taking the expectation of the update (Eq. 46)*

$$\mathbb{E}[\langle X \rangle_{t+1}] = \beta \mathbb{E}[\langle X \rangle_t] + (1 - \beta) \mathbb{E}[X_{t+1}]. \quad (51)$$

*Substituting Eq. (48),*

$$\mathbb{E}[\langle X \rangle_{t+1}] = \beta \mathbb{E}[\langle X \rangle_t] + (1 - \beta)\mu. \quad (52)$$

*To make the notation slightly “lighter”, use  $a_t = \mathbb{E}[\langle X \rangle_t]$ ,*

$$a_{t+1} = \beta a_t + (1 - \beta)\mu. \quad (53)$$

*We also have,*

$$a_t = \beta a_{t-1} + (1 - \beta)\mu. \quad (54)$$

$$a_{t-1} = \beta a_{t-2} + (1 - \beta)\mu. \quad (55)$$

*Substituting  $a_t$  into  $a_{t+1}$ ,*

$$a_{t+1} = \beta(\beta a_{t-1} + (1 - \beta)\mu) + (1 - \beta)\mu. \quad (56)$$

$$a_{t+1} = \beta^2 a_{t-1} + (1 - \beta)(1 + \beta)\mu \quad (57)$$

*Substituting  $a_{t-1}$ ,*

$$a_{t+1} = \beta^2(\beta a_{t-2} + (1 - \beta)\mu) + (1 - \beta)\mu(1 + \beta) \quad (58)$$

$$a_{t+1} = \beta^3 a_{t-2} + (1 - \beta)\mu(1 + \beta + \beta^2). \quad (59)$$

*Repeating until we get to  $a_0$ ,*

$$a_{t+1} = \beta^{t+1} a_0 + (1 - \beta)\mu \sum_{\tau=0}^t \beta^\tau. \quad (60)$$

*As we initialize at  $a_0 = 0$ ,*

$$a_{t+1} = (1 - \beta)\mu \sum_{\tau=0}^t \beta^\tau. \quad (61)$$

*Now, the  $\sum_{\tau=0}^t \beta^\tau$  term is a geometric series. You can look geometric series up. Or we can find the value ourselves. There’s a “trick” to the proof for a geometric series, which is to consider,*

$$(1 - \beta) \sum_{\tau=0}^t \beta^\tau = \sum_{\tau=0}^t \beta^\tau - \beta \sum_{\tau=0}^t \beta^\tau \quad (62)$$

*Now we can absorb the  $\beta$  into the sum,*

$$(1 - \beta) \sum_{\tau=0}^t \beta^\tau = \sum_{\tau=0}^t \beta^\tau - \sum_{\tau=1}^{t+1} \beta^\tau \quad (63)$$

Now, notice that almost all the terms cancel, except those for  $\tau = 0$  and  $\tau = T + 1$ ,

$$(1 - \beta) \sum_{\tau=0}^t \beta^\tau = 1 - \beta^{t+1}. \quad (64)$$

Substituting this back into our form for  $a_{t+1}$  gives,

$$a_{t+1} = (1 - \beta^{t+1})\mu \quad (65)$$

Of course, assuming that  $0 < \beta < 1$ , then as  $t + 1$  goes to infinity, we have,

$$\lim_{t \rightarrow \infty} \beta^{t+1} = 0 \quad \lim_{t \rightarrow \infty} a_{t+1} = \mu \quad (66)$$

The problem is what happens for smaller numbers of timesteps.

Part 2: Therefore, we can choose,

$$c_{t+1} = \frac{1}{1 - \beta^{t+1}} \quad (67)$$

Now, the expectation of our debiased estimator,  $d_{t+1}$  becomes,

$$\mathbb{E}[d_{t+1}] = c_{t+1} \mathbb{E}[\langle X \rangle_{t+1}] \quad (68)$$

Substituting  $a_{t+1} = \mathbb{E}[\langle X \rangle_{t+1}]$

$$\mathbb{E}[d_{t+1}] = c_{t+1} a_{t+1} \quad (69)$$

Substituting for  $c_{t+1}$  from Eq. (67) and  $a_{t+1}$  from Eq. (65),

$$\mathbb{E}[d_{t+1}] = \frac{1}{1 - \beta^{t+1}} (1 - \beta^{t+1})\mu = \mu. \quad (70)$$

So  $d_{t+1}$  has the correct expectation.