

Optimization Algorithms 2: RMSprop, Adadelata & Adam

ASTR 502
Maria Gabriela Cota Moreira
March 21st

Recap 1

AdaGrad

What it does?

Includes an adaptive learning rate that decreases the step-size with time

1. $\mathbf{g}_t = \partial_w \mathcal{L}(\mathbf{w}_t)$

Loss function gradient

2. $\mathbf{s}_t = \mathbf{s}_{t-1} + \mathbf{g}_t^2$

Learning rate scaling

3. $\mathbf{w}_t = \mathbf{w}_{t-1} - \frac{\eta}{\sqrt{\mathbf{s}_t + \epsilon}} \mathbf{g}_t$

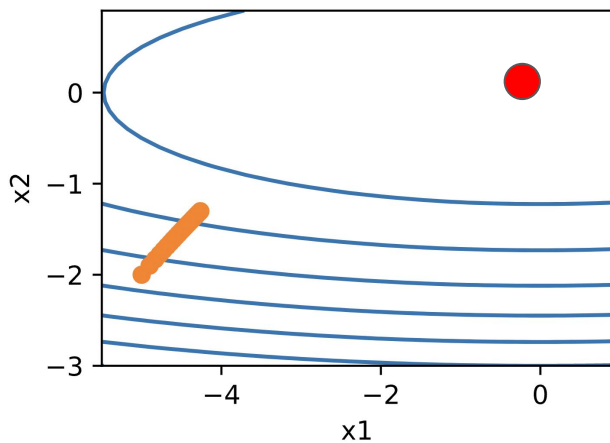
Parameter update

Issue: The more you update the less you update

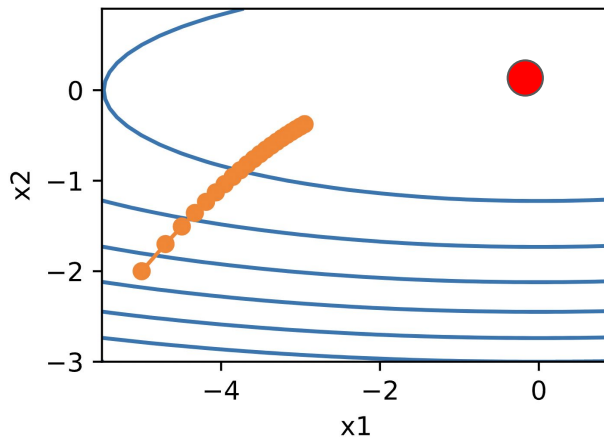
$$\mathbf{s}_t = \sum_{\tau=0}^t \mathbf{g}_{\tau}^2$$

$$\mathbf{w}_t = \mathbf{w}_{t-1} - \frac{\eta}{\sqrt{\mathbf{s}_t + \epsilon}} \mathbf{g}_t$$

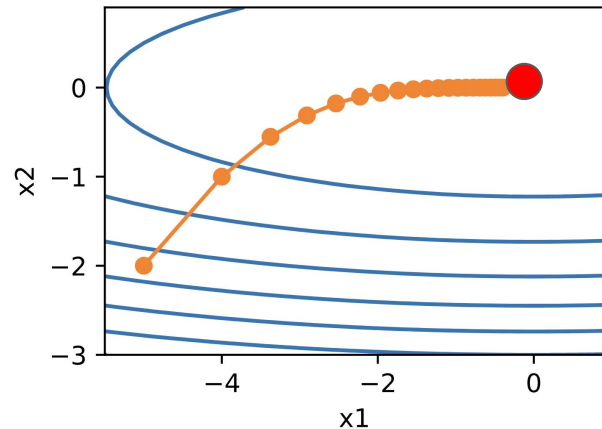
Slow and sensitive to the learning rate!



$\eta = 0.1$



$\eta = 0.3$



$\eta = 1.0$

RMSprop

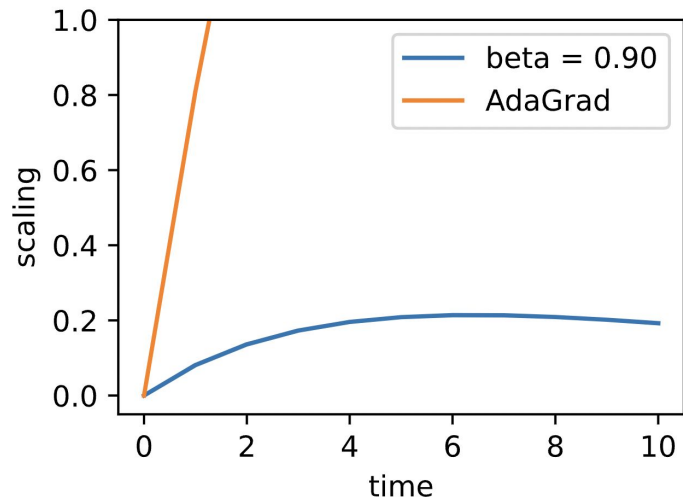
Root-Mean-Square Propagation

What's new?

Includes a damping parameter to AdaGrad, β w/ same $\mathbf{w}_t = \mathbf{w}_{t-1} - \frac{\eta}{\sqrt{\mathbf{s}_t + \epsilon}} \mathbf{g}_t$

Leaky average of the squared gradients

$$\mathbf{s}_t = \beta \mathbf{s}_{t-1} + (1 - \beta) \mathbf{g}_t^2$$



Keeps the squares under a manageable size the whole time !

Initial values
and params

$$\mathbf{s}_0 = 0$$

$$\beta = 0.9$$

$$\epsilon = 10^{-8}$$

How it works? We give more importance to recent gradients

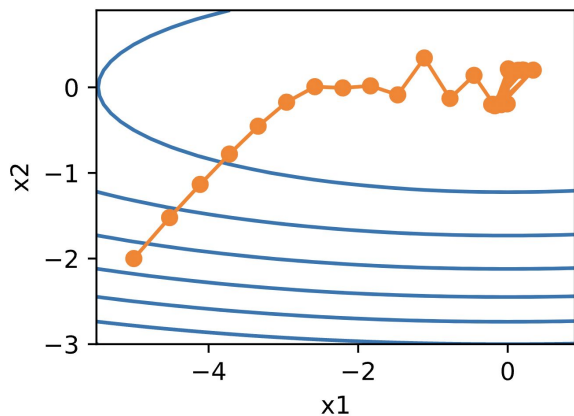
$$0 < \beta < 1$$

$$\beta = 0.9$$

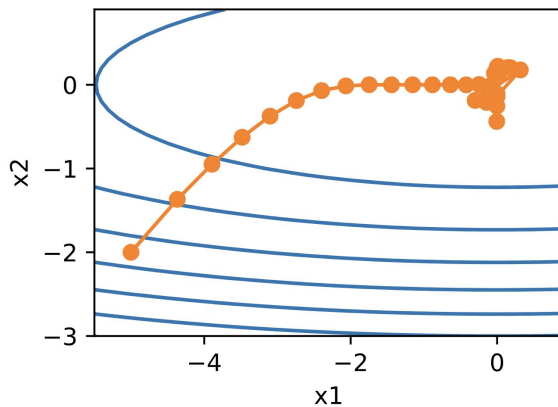
$$\mathbf{s}_t = \beta \mathbf{s}_{t-1} + (1 - \beta) \mathbf{g}_t^2$$

$$\mathbf{s}_t = (1 - \beta)(\mathbf{g}_t^2 + \beta \mathbf{g}_{t-1}^2 + \beta^2 \mathbf{g}_{t-2}^2 + \beta^3 \mathbf{g}_{t-3}^2 + \dots)$$

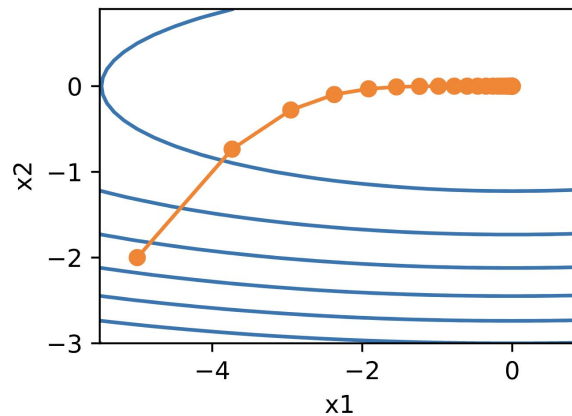
lower weight to previous gradients
variance



$$\beta = 0.3$$



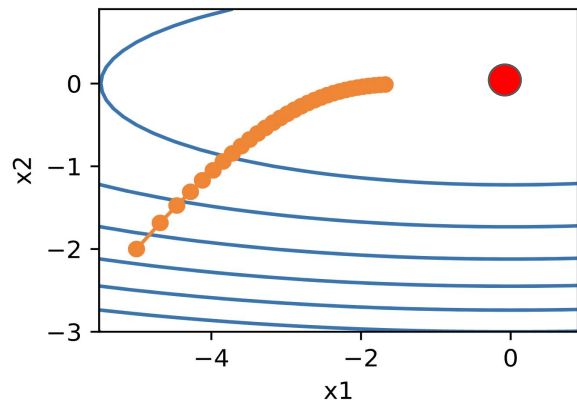
$$\beta = 0.6$$



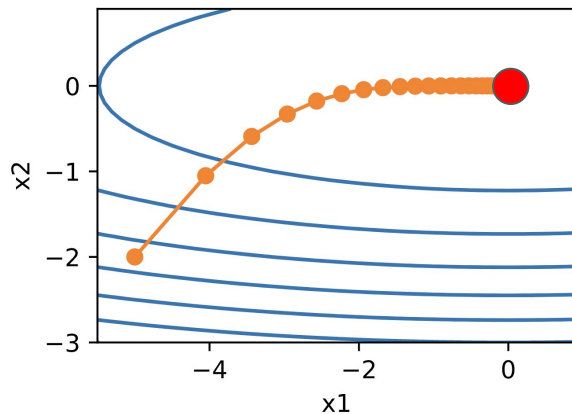
$$\beta = 0.9$$

Potential issue

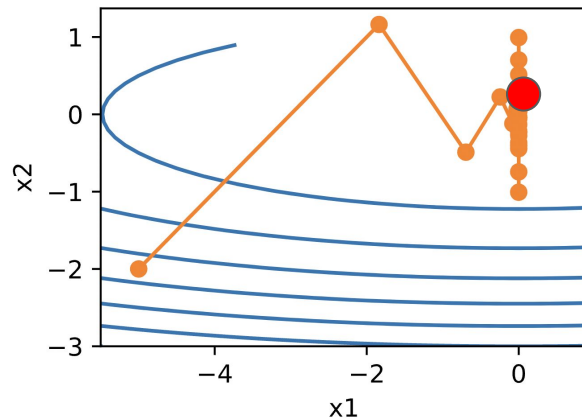
Learning rate is still manual, because the suggested value is not always appropriate for every task.



$$\eta = 0.1$$



$$\eta = 0.3$$



$$\eta = 1.0$$

ADADELTA

An Adaptive Learning Rate Method

What's new?

I) Also Includes a damping parameter to AdaGrad

II) Address the need for a manually selected global learning rate (kinda...)

1. $\mathbf{g}_t = \partial_w \mathcal{L}(\mathbf{w}_t)$

2. $\mathbf{s}_t = \beta \mathbf{s}_{t-1} + (1 - \beta) \mathbf{g}_t^2$



Same as RMSprop

3. $\mathbf{g}'_t = \frac{\sqrt{\Delta \mathbf{w}_{t-1} + \epsilon}}{\sqrt{\mathbf{s}_t + \epsilon}} \odot \mathbf{g}_t$



Step-size

4. $\Delta \mathbf{w}_t = \beta \Delta \mathbf{w}_{t-1} + (1 - \beta) \mathbf{g}'_t{}^2$



Leaky average of the step-size

5. $\mathbf{w}_t = \mathbf{w}_{t-1} - \mathbf{g}'_t$



Parameter update

Initial values

$$\Delta \mathbf{w}_0 = 0$$

$$\mathbf{s}_0 = 0$$

$$\beta = 0.9$$

Motivation/How it works

Numerator arises from unit correction.

In AdaGrad,

$$\mathbf{g}'_t = -\frac{\eta}{\sqrt{\mathbf{s}_t + \epsilon}} \odot \mathbf{g}_t \quad \text{is unitless} \quad \mathbf{g}'_t \propto \frac{\frac{\partial \mathcal{L}}{\partial \mathbf{w}}}{\sqrt{\left(\frac{\partial \mathcal{L}}{\partial \mathbf{w}}\right)^2}}$$

But it should have units of \mathbf{w} since $\mathbf{w}_t = \mathbf{w}_{t-1} - \mathbf{g}'_t$. Then,

$$\mathbf{g}'_t \propto \frac{\boxed{\Delta \mathbf{w}_t} \frac{\partial \mathcal{L}}{\partial \mathbf{w}}}{\sqrt{\left(\frac{\partial \mathcal{L}}{\partial \mathbf{w}}\right)^2}} \propto \text{units of } \mathbf{w}_t$$

$$\mathbf{w}_t = \mathbf{w}_{t-1} - \frac{\boxed{\sqrt{\Delta \mathbf{w}_{t-1} + \epsilon}}}{\sqrt{\mathbf{s}_t + \epsilon}} \odot \mathbf{g}_t$$

It acts as a “momentum” term

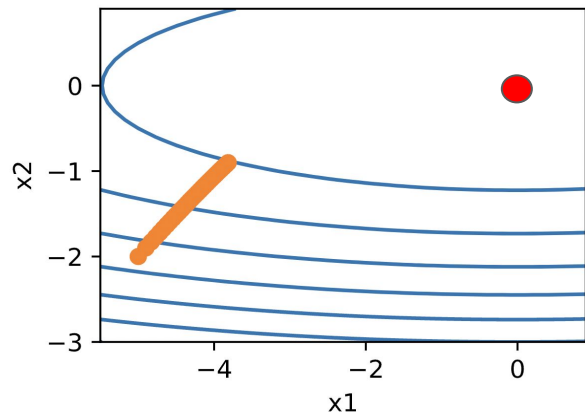
Potential issues (that I ran into)

When I chose a small ϵ the algorithm converges slowly

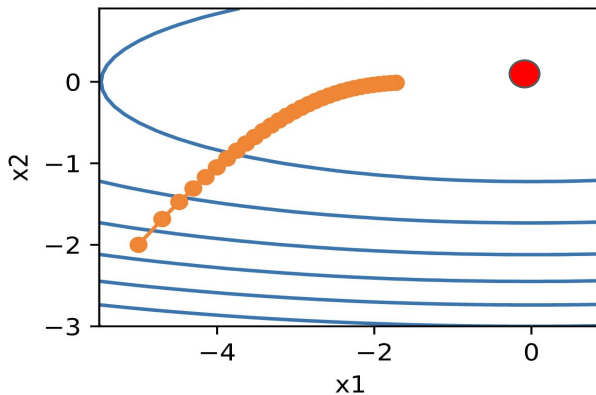
$$\mathbf{g}'_t = \frac{\sqrt{\Delta \mathbf{w}_{t-1} + \epsilon}}{\sqrt{\mathbf{s}_t + \epsilon}} \odot \mathbf{g}_t$$

Is it really “learning-rate free”?

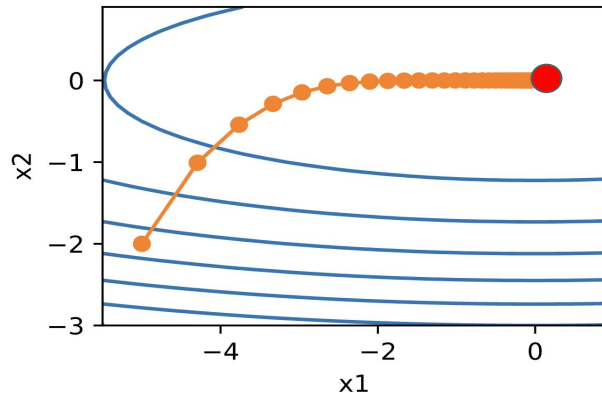
Changing the initial $\Delta \mathbf{w}_0$ also leads to a similar effect: <https://akyriillidis.github.io/notes/AdaDelta>



$\epsilon = 0.001$



$\epsilon = 0.01$



$\epsilon = 0.1$

Potential issues (*that I ran into*)

ϵ is behaving as a **learning rate**. But why?

$$\mathbf{g}'_t = \frac{\sqrt{\Delta \mathbf{w}_{t-1} + \epsilon}}{\sqrt{\mathbf{s}_t + \epsilon}} \odot \mathbf{g}_t$$

at the first
iteration/epoch (t=1)
you have

$$\mathbf{s}_1 = (1 - \beta) \mathbf{g}_1^2$$

$$\Delta \mathbf{w}_0 = 0$$

then,

$$\mathbf{g}'_1 = \frac{\sqrt{\epsilon}}{\sqrt{\mathbf{s}_1 + \epsilon}} \odot \mathbf{g}_1 \longrightarrow \Delta \mathbf{w}_1 = (1 - \beta) \mathbf{g}_1'^2$$

$\nearrow \eta = \sqrt{\epsilon}$

Recap 2

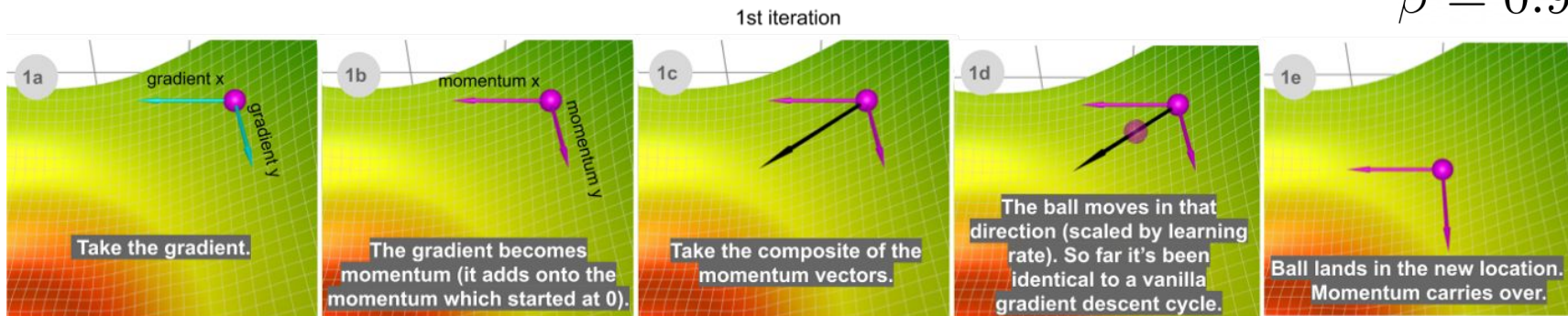
Momentum

Momentum accelerates motion towards the minima

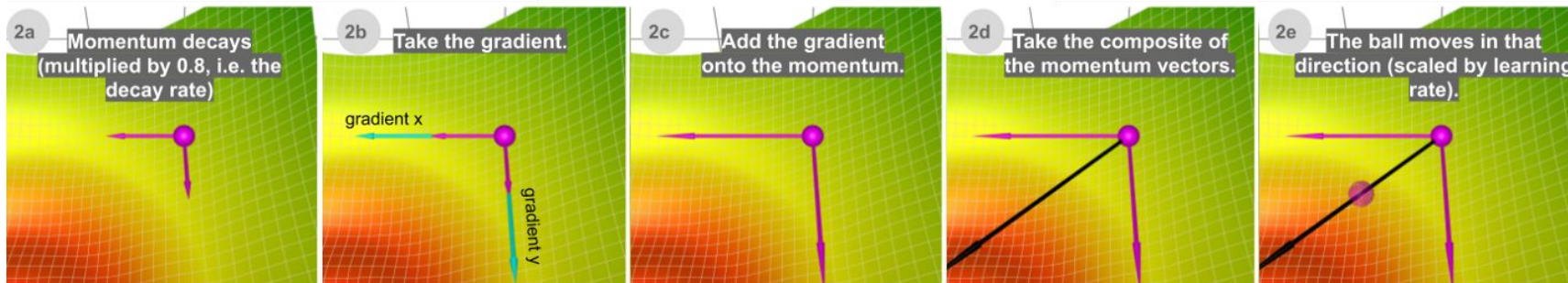
$$\mathbf{v}_t = \beta \mathbf{v}_{t-1} + (1 - \beta) \mathbf{g}_t$$

$$\mathbf{w}_t = \mathbf{w}_{t-1} + \eta \mathbf{v}_t$$

$$\beta = 0.9$$



2nd iteration (a typical momentum descent cycle)



3rd iteration starts, carrying over the momentum, so on and so forth...

Adam

Adaptive Moment Estimation

What's new?

Utilize the **momentum** concept and **adaptive learning rate** from RMSprop

Momentum

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

$$\mathbf{w}_t = \mathbf{w}_{t-1} + \eta \mathbf{v}_t$$



RMSprop

$$\mathbf{s}_t = \beta_2 \mathbf{s}_{t-1} + (1 - \beta_2) \mathbf{g}_t^2$$

$$\mathbf{w}_t = \mathbf{w}_{t-1} - \frac{\eta}{\sqrt{\mathbf{s}_t + \epsilon}} \mathbf{g}_t$$

bias correction



Adam

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

$$\mathbf{w}_t = \mathbf{w}_{t-1} - \frac{\eta}{\sqrt{\hat{\mathbf{s}}_t + \epsilon}} \odot \hat{\mathbf{v}}_t$$

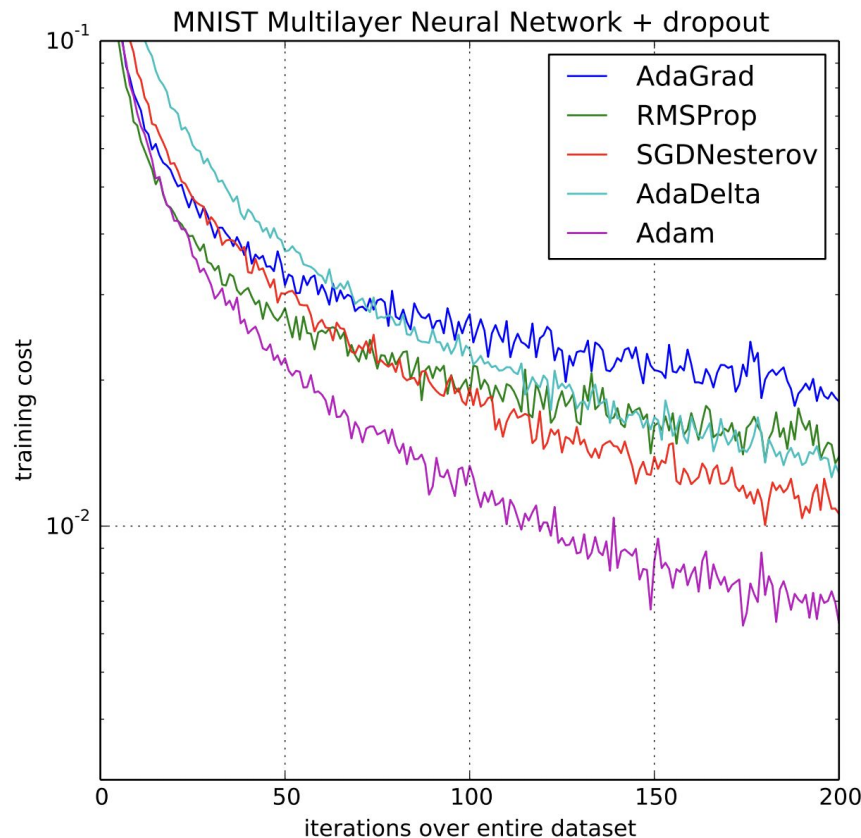
$$\beta_1 = 0.9 \quad \beta_2 = 0.999 \quad \epsilon = 10^{-8}$$

How it works?

Momentum accelerates motion

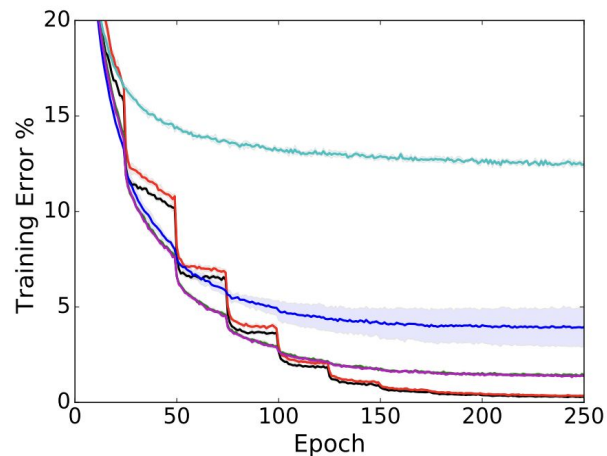
RMSprop controls oscillations

Taking advantage of both is
what makes Adam powerful and
popular

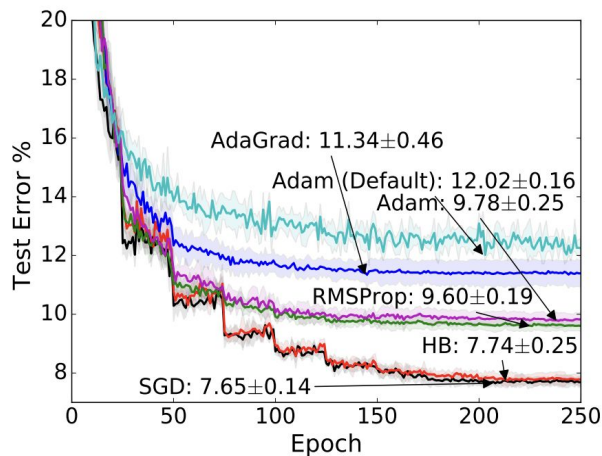


Potential issues Don't trust on the default settings too much

Many hyperparameters to tune: $\beta_1, \beta_2, \eta, \epsilon, \mathbf{s}_0$ and \mathbf{v}_0



(a) CIFAR-10 (Train)



(b) CIFAR-10 (Test)

The authors observe that the solutions found by adaptive methods generalize worse than SGD, even when these solutions have better training performance

Figure 1: Training (left) and top-1 test error (right) on CIFAR-10. The annotations indicate where the best performance is attained for each method. The shading represents \pm one standard deviation computed across five runs from random initial starting points. In all cases, adaptive methods are performing worse on both train and test than non-adaptive methods.

<https://arxiv.org/pdf/1705.08292.pdf>

References

Dive into Deep Learning, Ch.11.8-11.10 (<https://d2l.ai/>)

<https://akyrillidis.github.io/notes/AdaDelta>

M. D. Zeiler “ADADELTA: An adaptive learning rate method” ([arXiv:1212.5701](https://arxiv.org/abs/1212.5701), 2012)

D. P. Kingman and J. Lei Ba*“Adam: A method for stochastic optimization” ([arXiv:1412.6980](https://arxiv.org/abs/1412.6980))

<https://towardsdatascience.com/a-visual-explanation-of-gradient-descent-methods-momentum-adagrad-rmsprop-adam-f898b102325c>

Gradient Descent Visualization (https://github.com/lilipads/gradient_descent_viz)

A. C. Wilson et al. “The Marginal Value of Adaptive Gradient Methods in Machine Learning” ([arXiv:1705.08292](https://arxiv.org/abs/1705.08292))