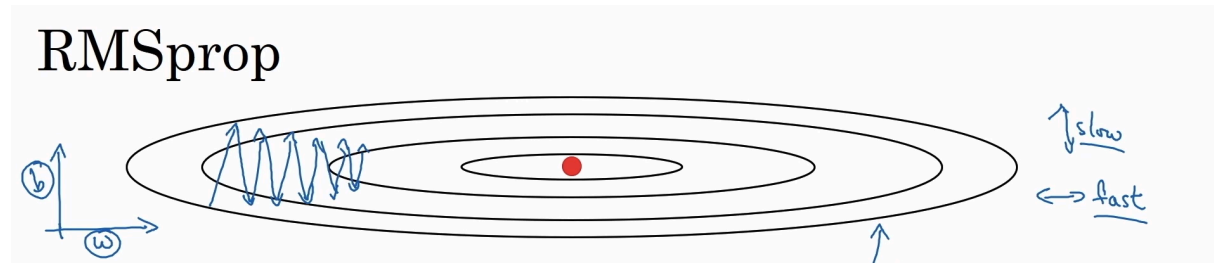


# RMSProp

RMSProp → Root Mean Square Propagation



Want to:

Slow down learning in the b direction and increase learning in the W direction.

On iteration  $t$ :

Compute  $dW, db$  on current mini-batch

$S_{dW} = \beta S_{dW} + (1-\beta) \frac{dW^2}{\text{element-wise}}$

$S_{db} = \beta S_{db} + (1-\beta) db^2$

$W := W - \alpha \frac{dW}{\sqrt{S_{dW}}}$

$b := b - \alpha \frac{db}{\sqrt{S_{db}}}$

Do note that the squaring is done element wise!

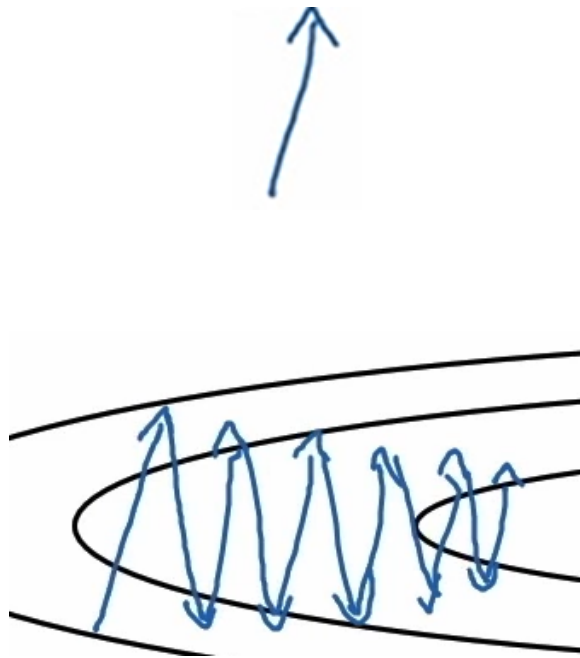
So what this is doing is keeping an EWA of the SQUARES of the derivatives.

Then look carefully at the update rule! We don't directly update by the EWA!

So, goal is to increase learning in the W direction and decrease/dampen learning in the b direction.

So what we hope is that  $S_{dW}$  is small so that when you divide by a smaller number, then the update to W is large. And we hope  $S_{db}$  is large so that when you divide by a larger number, then the update to b is small in order to slow down updates in the vertical direction.

And indeed, the  $dW$ 's are much smaller than the  $db$ 's!



Here,  $db$  is much larger compared to  $dW$

Therefore oscillations in the vertical direction get damped out!

Think of this as adaptive learning rates!