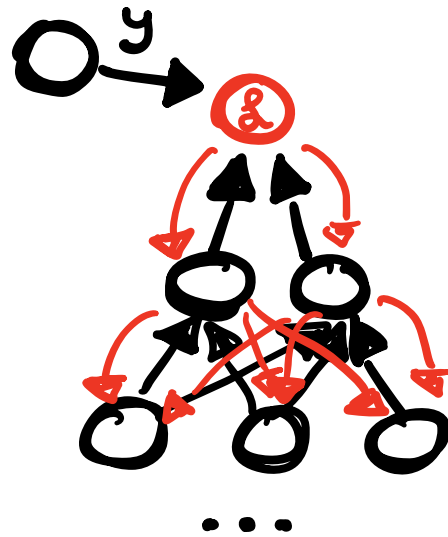


DS 4440

Optimizer MATTERS



Backprop yields gradients for all parameters efficiently ($\nabla_W L$).

Cool! ... Why did we want these again? For Gradient Descent!

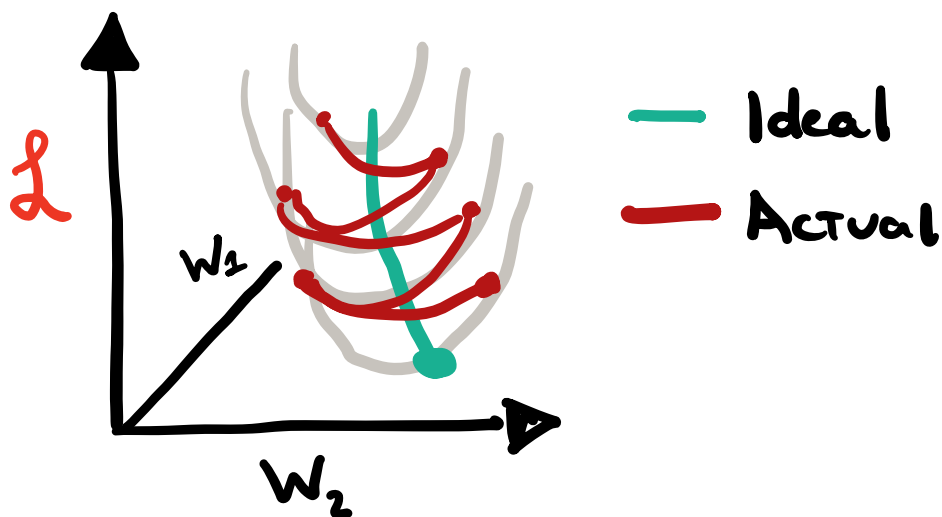
$$\hat{W}_{\tau+1} \leftarrow \hat{W}_{\tau} - \underset{\substack{\text{learning} \\ \text{rate}}}{\alpha} \nabla_W L(\hat{W}_{\tau}, x, y)$$

Variants

Stochastic 1 instance (x, y) per step

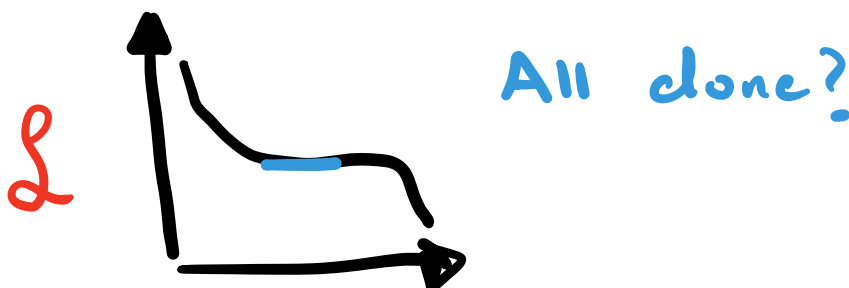
Batch Sets of (x, y) ; $\frac{1}{|B|} \sum \nabla_w \mathcal{L}$

Basic Gradient Descent is Myopic
in that it zig zags \rightarrow Slow Convergence



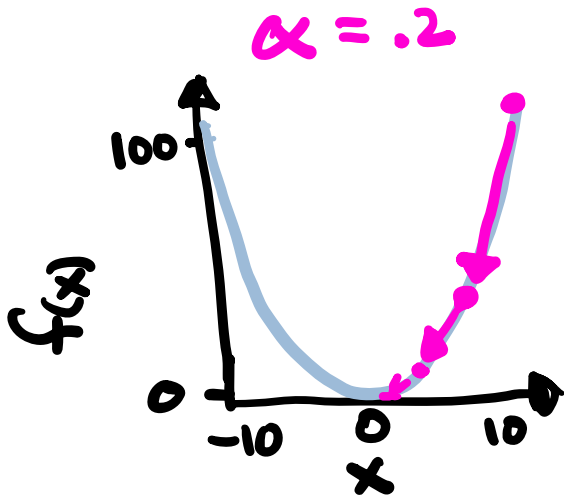
Normalizing inputs can help a bit.

A general challenge: Saddle points



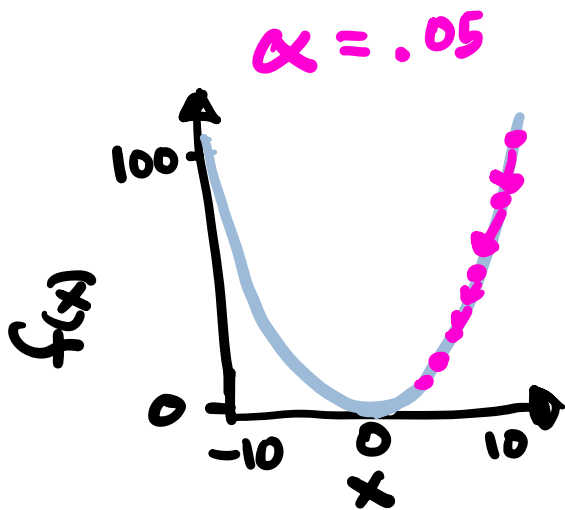
The Importance of learning rate α

$$f(x) = x^2 \quad \frac{d}{dx} f(x) = 2x \quad x_{\text{INIT}} = 10$$



$$\begin{aligned} .2 \cdot (2 \cdot 10) &= 4 & x &= 6 \\ .2 \cdot (2 \cdot 6) &= 2.4 & x &= 3.6 \\ .2 \cdot (2 \cdot 2.4) &= .96 & x &\approx 2.6 \\ &\dots & & \end{aligned}$$

BUT $\alpha = .05$



$$\begin{aligned} .05 \cdot (2 \cdot 10) &= 1 & x &= 9 \\ .05 \cdot (2 \cdot 9) &= .9 & x &= 8.1 \\ .05 \cdot (2 \cdot 8.1) &= .81 & x &\approx 7.3 \\ &\dots & & \text{Slloooooow} \end{aligned}$$

MOMENTUM

BUT adjust locally (e.g. .1)

$$v_{\tau} \leftarrow \gamma v_{\tau-1} - \alpha \nabla_w L(\hat{w}_{\tau}, x, y)$$

Keep going in direction
you've been going ($\gamma = .9$, e.g.)

$$\hat{W}_{\tau+1} \leftarrow \hat{W}_{\tau} + v_{\tau}$$

We can even anticipate where we are going:

$$v_{\tau} \leftarrow \gamma v_{\tau-1} - \alpha \mathcal{L}(\underbrace{\hat{W}_{\tau} + \gamma v_{\tau-1}}_{\text{Move toward accumulated } \nabla_{\mathbf{w}}}, x, y)$$

This is Nesterov Momentum.

(See examples in CoLab.)

So far we have assumed a single learning rate (α) for all parameters.

$$\text{Let } g_{\tau,j} \equiv \nabla_{\mathbf{w}} \mathcal{L}(\hat{W}_{\tau})_j$$

$$\hat{W}_{(\tau+1),j} \leftarrow \hat{W}_{\tau,j} - \alpha_{\tau,j} g_{\tau,j} \quad | \quad \text{"AdaGrad"}$$

But how to set $\alpha_{\tau,j}$?

INTUITION Move slower for params

We have updated a bunch

$$\alpha_{\tau,j} \doteq \frac{\alpha}{\sqrt{G_{\tau,j}} + \epsilon} \quad \text{--- Avoids div. by } \emptyset.$$

Sum of grads for j .

$$\sum_{\tau} \left[\nabla_w \mathcal{L}(\hat{W}_{\tau}, x, y) \right]_j^2$$

Cumulative updates slow rate.

RMSProp Extends AdaGrad

by keeping a decaying average of the squares of ∇ terms.

$$S_{\tau,j} \doteq \lambda S_{(\tau-1),j} + (1-\lambda) g_{\tau,j}^2$$

Then we divide α by $\sqrt{S_{\tau,j}}$.

Adam is maybe the go-to optimizer now. We keep two vars per param:

$$m_{\tau} \leftarrow \beta_2 m_{\tau-1} + (1-\beta_2) g_{\tau} \approx \text{MOMENTUM}$$

$$v_{\tau} \leftarrow \beta_2 v_{\tau-1} + (1-\beta_2) g_{\tau}^2 \approx \text{RMS Prop}$$

Then

$$\hat{w}_{(\tau+1),j} \leftarrow \hat{w}_{\tau,j} - \frac{\alpha}{\sqrt{v_{\tau,j}} + \epsilon} \cdot m_{\tau,j}$$

Slightly Simplified; ignores bias-correction.

Learning rate decay

Loss Surface

Big Steps TO
START \rightarrow Smaller
over Time



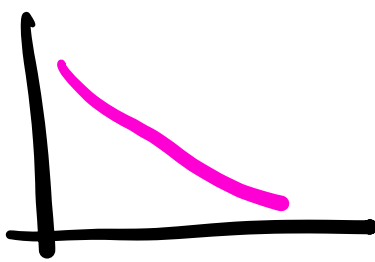
Assume noisy
minibatches

Big α might
Jump too much

Shrinking α at each
Step might help

$$\alpha = \frac{1}{1 + d \cdot \text{epoch}} \cdot \alpha_0$$

decay rate

if $d=1 \rightarrow$ 

Can also do exponential decay

$$\alpha = (0.9)^{\text{epoch}} \cdot \alpha_0$$

Note on hyperparameter Tuning

- Choose random values (not grid)