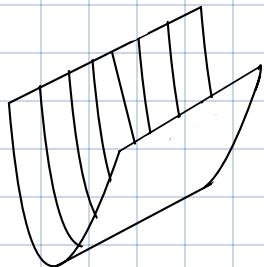
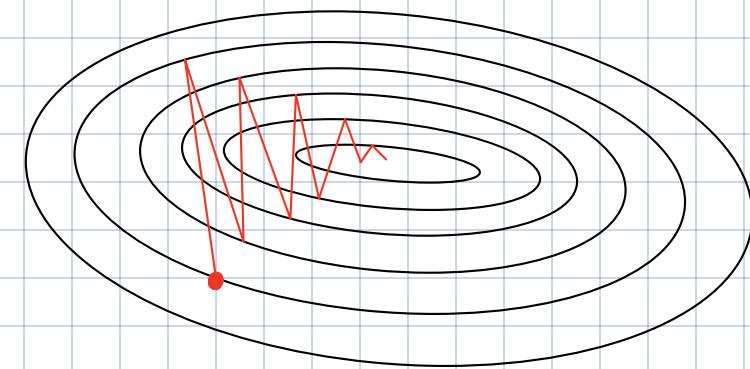


# Fancier Optimizations:



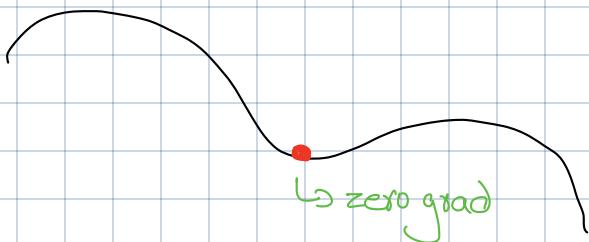
- Sensitive in vertical direction & not in horizontal

SHD:

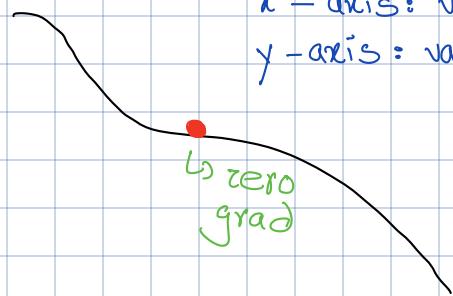


↳ Direction of gradient doesn't align with direction towards minima

Local minima / Saddle points



$x$ -axis: value of a param  
 $y$ -axis: value of loss



↳ SHD will get stuck

Saddle point:



A point where loss increases in one direction & decreases in another

local minima:



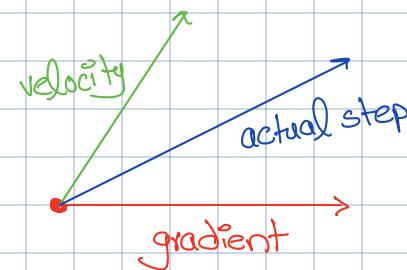
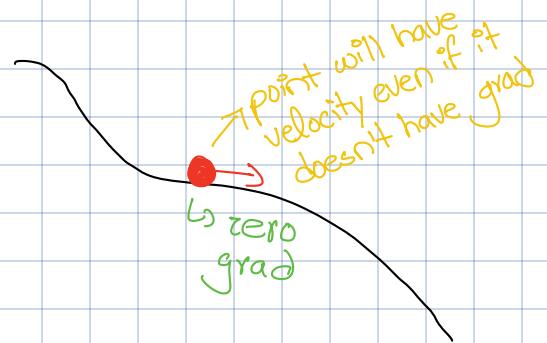
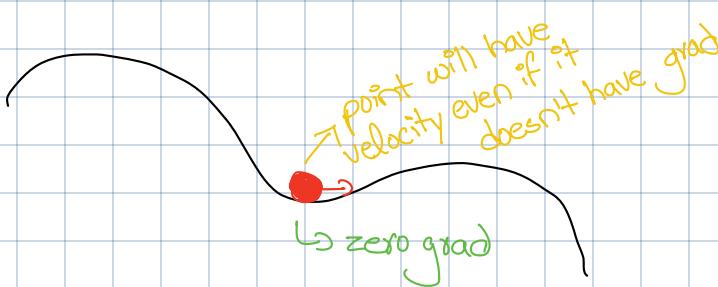
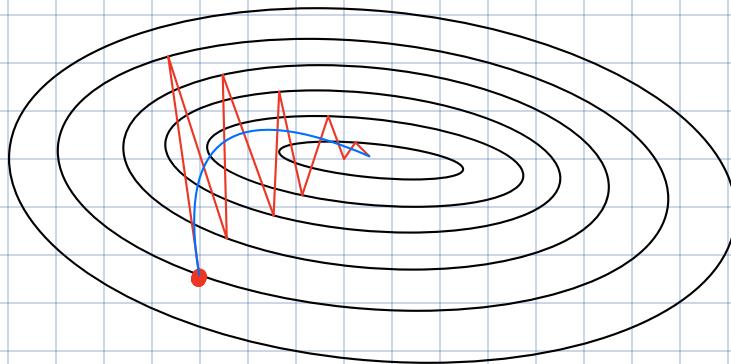
Anywhere else the loss increases

SGD + momentum:

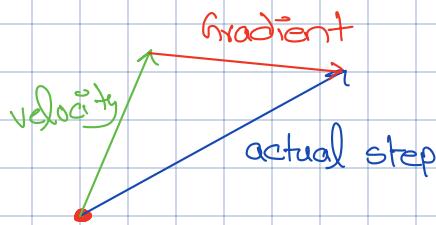
friction - 0.9 or 0.99

$$v_{t+1} = \gamma v_t + \nabla f(x_t)$$

$$x_{t+1} = x_t - \alpha v_{t+1}$$



Nesterov Momentum:



↳ Step in the point where velocity will take you

↳ Calculate the grad at that point

↳ Come back to your original point & mix together those two

$$v_{t+1} = \gamma v_t - \alpha \nabla f(x_t + \gamma v_t)$$

$$x_{t+1} = x_t + v_{t+1}$$

OR

$$\text{Let } \bar{x}_t = x_t + \gamma v_t$$

$$v_{t+1} = \gamma v_t - \alpha \nabla f(\bar{x}_t)$$

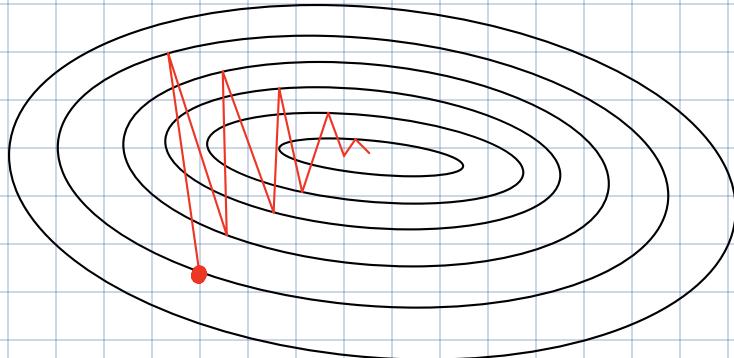
$$\begin{aligned} \bar{x}_{t+1} &= \bar{x}_t - \gamma v_t + (1-\gamma)v_{t+1} \\ &= \bar{x}_t + v_{t+1} + \gamma(v_{t+1} - v_t) \end{aligned}$$

Adagrad:

$$\text{grad\_sqj} = \text{grad\_sqj} + (\nabla f(x_t))^2$$

$$x_{t+1} = x_t - \alpha \nabla f(x_t)$$

$$\frac{\nabla \text{grad-sq} + \epsilon p}{\hookrightarrow 1e-7}$$



- ↳ In vertical direction, grads are high so we divide by the large gradient  $\downarrow$  then take small update step
- ↳ In horizontal direction, grads are small so we divide by the small gradient  $\downarrow$  then take a large update step
- ↳ One problem is that over time, grad-sq becomes so big that our update step becomes too small.

RMSprop:

$$\text{grad-sq} = \text{grad-sq} * \text{decay\_rate} + (1 - \text{decay\_rate}) * \nabla f(x_t)^2$$

$$x_{t+1} = x_t - \alpha \frac{\nabla f(x_t)}{\sqrt{\text{grad-sq} + \epsilon p}}$$

$$\hookrightarrow 1e-7$$

Adam:

- ↳ combines SGD + momentum  $\downarrow$  RMSprop

$$v_{t+1} = \beta_1 * v_t + \nabla f(x_t)$$

$$\text{grad-sq} = \text{grad-sq} * \beta_2 + (1 - \beta_2) * \nabla f(x_t)^2$$

$$x_{t+1} = x_t - \alpha \frac{v_{t+1}}{\sqrt{\text{grad-sq} + \epsilon p}}$$

↳ At  $t=0$ ,  $\nabla v \propto \text{grad-sq} = 0$ . So,  $\text{grad-sq}$  becomes  $\approx 0$  & we end up taking large steps at the beginning

So, we end up with:

$$v_{t+1} = \beta_1 v_t + \nabla f(x_t)$$

$$\text{grad-sq} = \text{grad-sq} * \beta_2 + (1 - \beta_2) * \nabla f(x_t)^2$$

$$v_{t+1\text{-correction}} = \frac{v_{t+1}}{(1 - \beta_1^t)}$$

$$\text{grad-sq-correction} = \frac{\text{grad-sq}}{(1 - \beta_2^t)}$$

$$x_{t+1} = x_t - \alpha \frac{v_{t+1\text{-correction}}}{\sqrt{\text{grad-sq-correction}} + \epsilon_p}$$

Learning rate:

↳ Decay  $\alpha$  over time:

• Exponential decay:

$$\alpha = \alpha_0 e^{-kt}$$

•  $1/t$  decay:

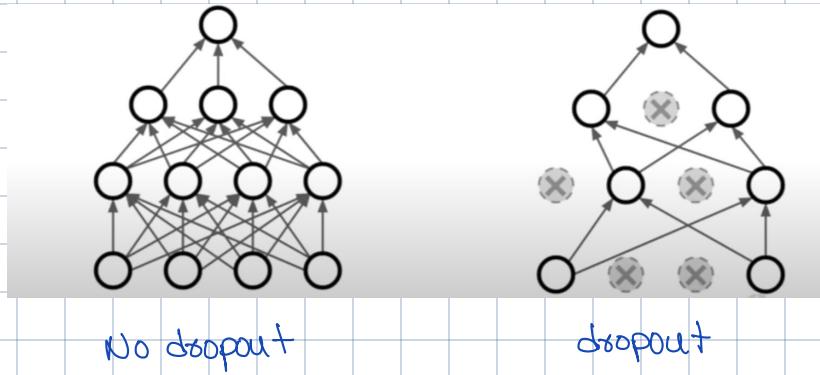
$$\alpha = \frac{\alpha_0}{(1+kt)}$$

• Step decay:

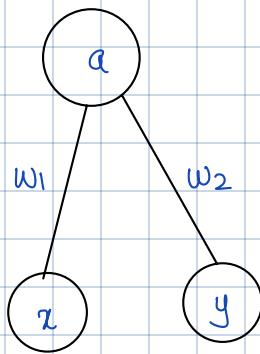
↳ Decay  $\alpha$  by  $\frac{1}{2}$  every few epochs

## Dropout:

↳ During forward pass, set randomly some neurons to 0



At test time:



$$E[a] = w_1x + w_2y$$

$$\begin{aligned} E[a] &= \frac{1}{4}(w_1x + w_2y) + \frac{1}{4}(w_1x + 0y) \\ &\quad + \frac{1}{4}(0x + 0y) + \frac{1}{4}(0x + w_2y) \\ &= \frac{1}{2}(w_1x + w_2y) \end{aligned}$$

$$z = (wx + b) * p$$

## DropConnect:

↳ Randomly zero out values of  $w$

## Fractional Pooling:

↳ Randomize pooling regions

## Stochastic depth:

↳ Randomly drop layers

Assignment helper:

