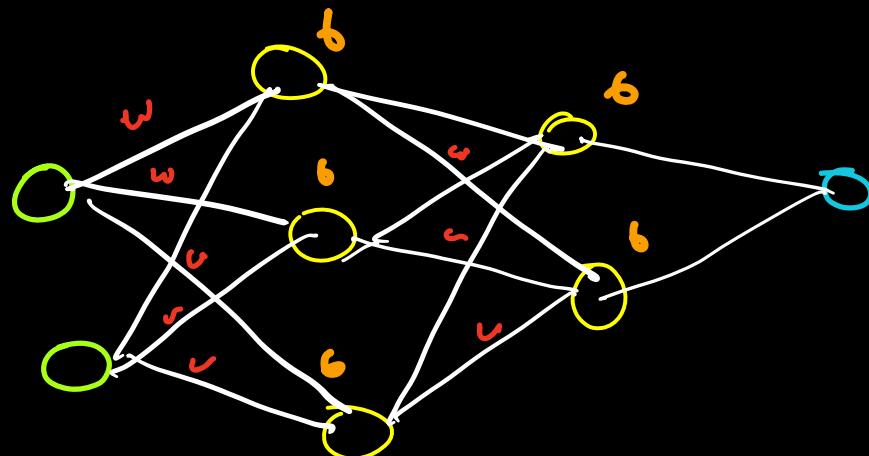


Weight Initialisation &

Optimisers

[Neural Network]

Q: why random init? Why not constant?



$\rightarrow 0$   
 $\rightarrow 1$   
 $\rightarrow 75\%$  ?

if all weights are same in beginning  
then, all gradients will be same for  
every-one. Then all updates will be  
same. Hence network will never learn.  
 $\rightarrow$  A-symmetry is very important.

## Random Init

$$w_{ij}^k \sim N(0, \sigma) \quad \text{OR}$$

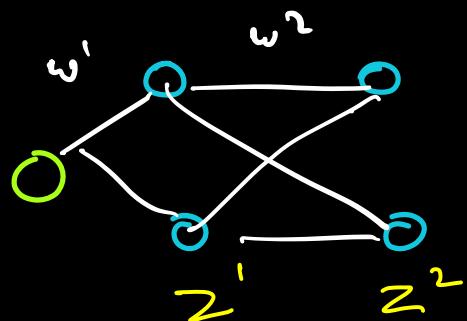
$$w_{ij}^k \sim U(-1, 1)$$

But this causes the problem of  
exploding gradients in deep NN.

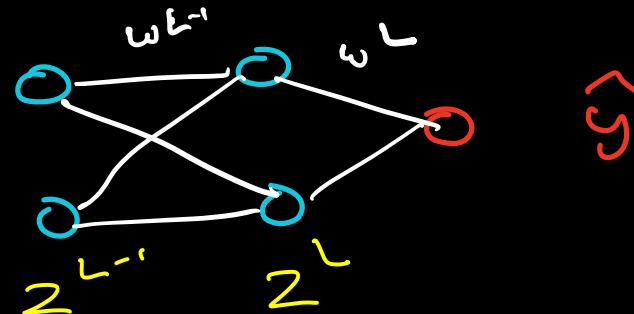
↓  
large number  
of hidden layers.

## Vanishing / Exploding Gradients

Let's take a deep NN. (simple one)



- - -



Let's say all  $b_i = 0$   
and linear activation, so  $g(z) = z$

$$\therefore z' = w' x$$

$$\therefore z^2 = w^2 z' \rightarrow w^2 w' x$$

$$\therefore y = w^L \cdot w^{L-1} \cdots w^2 \cdot \underbrace{w' x}_1$$

$$y = \sum_i w_i g(x_i) + b$$

$\therefore b = 0$   
 $g(2) = 2$

We see too many weight multipliers  
 So if  $w_i > 1$ , then,

Since  $L$  is large

$$(1.1)^L \rightarrow \infty \quad \text{Exploding}$$

$$(0.8)^L \rightarrow 0 \quad \text{Vanishing}$$

This will also happen in back prop!!

## Solution

1. ReLU  $\rightarrow$  Since its gradient is '1' on +ve side, it mitigates.  $1^n = 1 \checkmark$

However since it does not dampen activation, if weights are large it can cause exploding gradient.

2. Smart weight initialisation:

There is not a strong theoretical reason but the following have been shown to work experimentally

## 1. He Init

$$\rightarrow N(0, \frac{2}{n^{e-1}})$$

$\hookrightarrow$  # neurons in  
prev layer

$$\rightarrow U\left[\frac{-\sqrt{6}}{\sqrt{n^{e-1}}}, \frac{\sqrt{6}}{\sqrt{n^{e-1}}}\right]$$

## 2. Global Init.

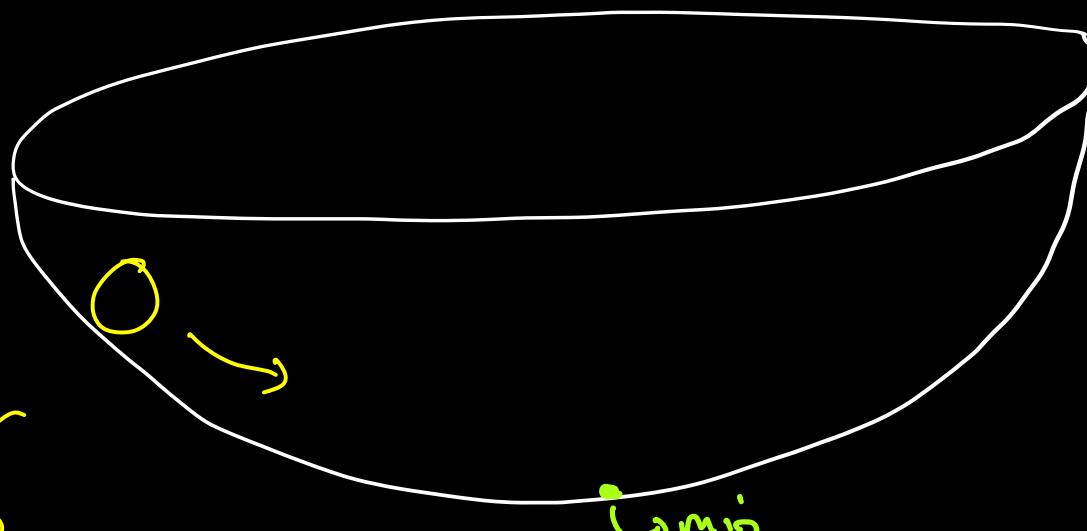
$$\rightarrow N(0, \frac{2}{n^{e-1} + n^e})$$

$$\rightarrow U\left(\frac{-\sqrt{6}}{\sqrt{n^{e-1} + n^e}}, \frac{\sqrt{6}}{\sqrt{n^{e-1} + n^e}}\right)$$

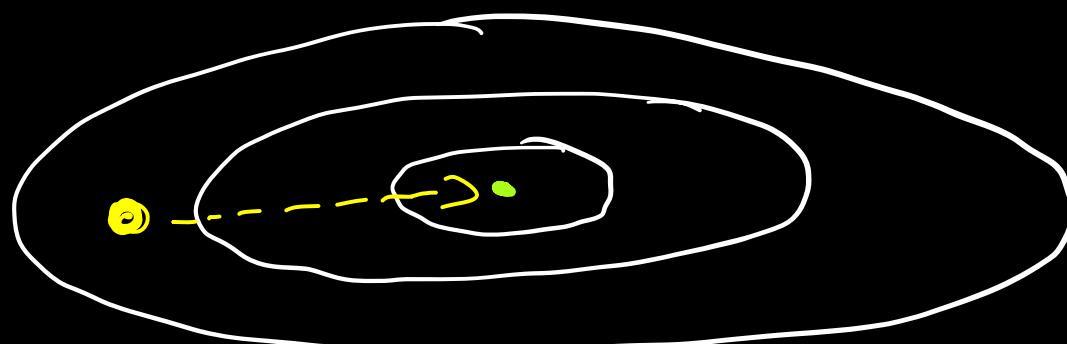
$\rightarrow$  code

# Optimizers

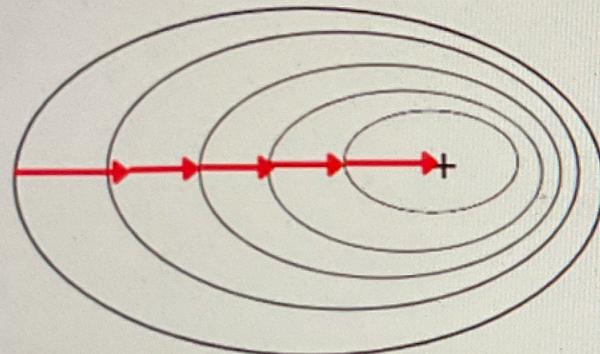
→ imagine a bowl



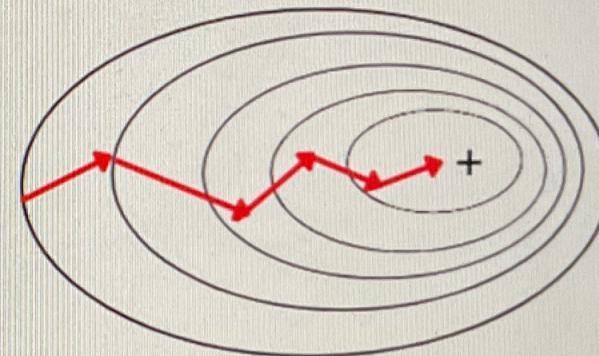
contour  
plot  
(shadow)



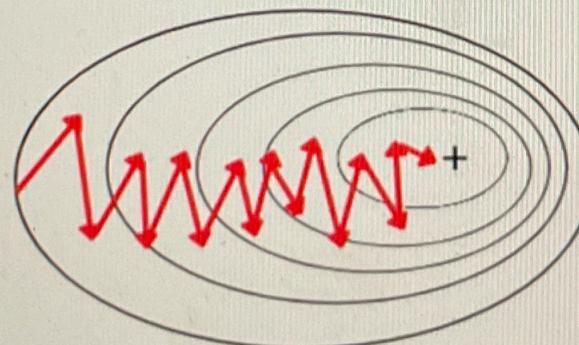
Batch Gradient Descent

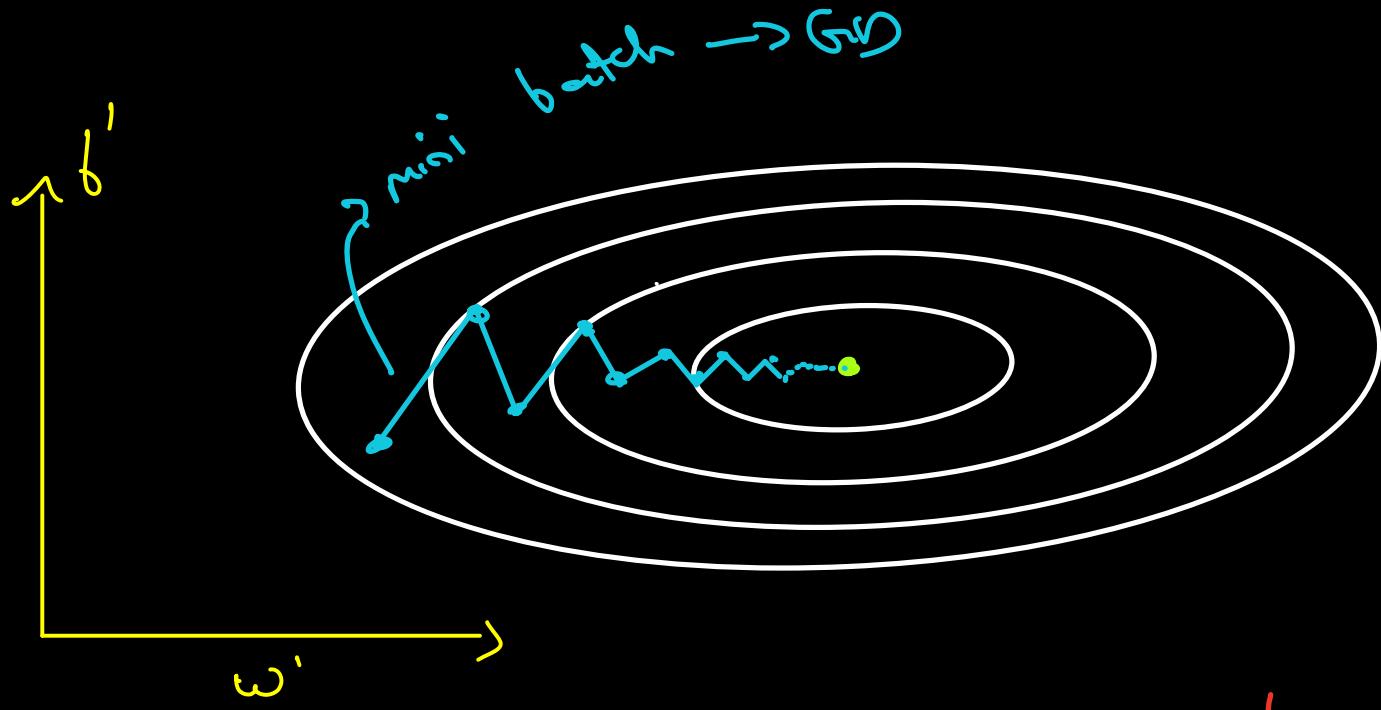


Mini-Batch Gradient Descent



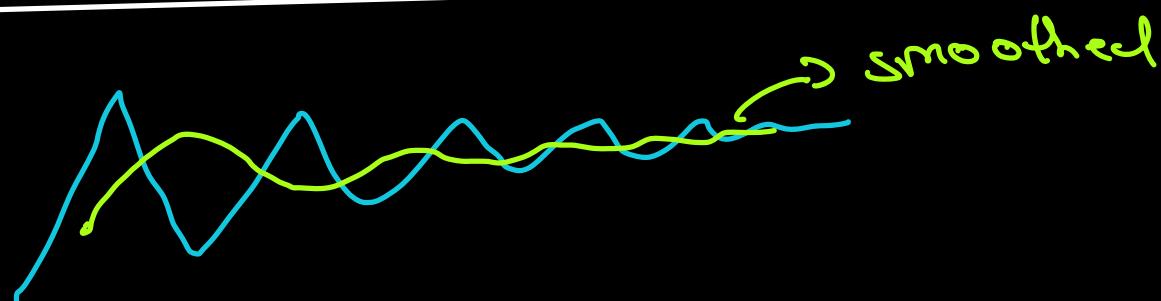
Stochastic Gradient Descent





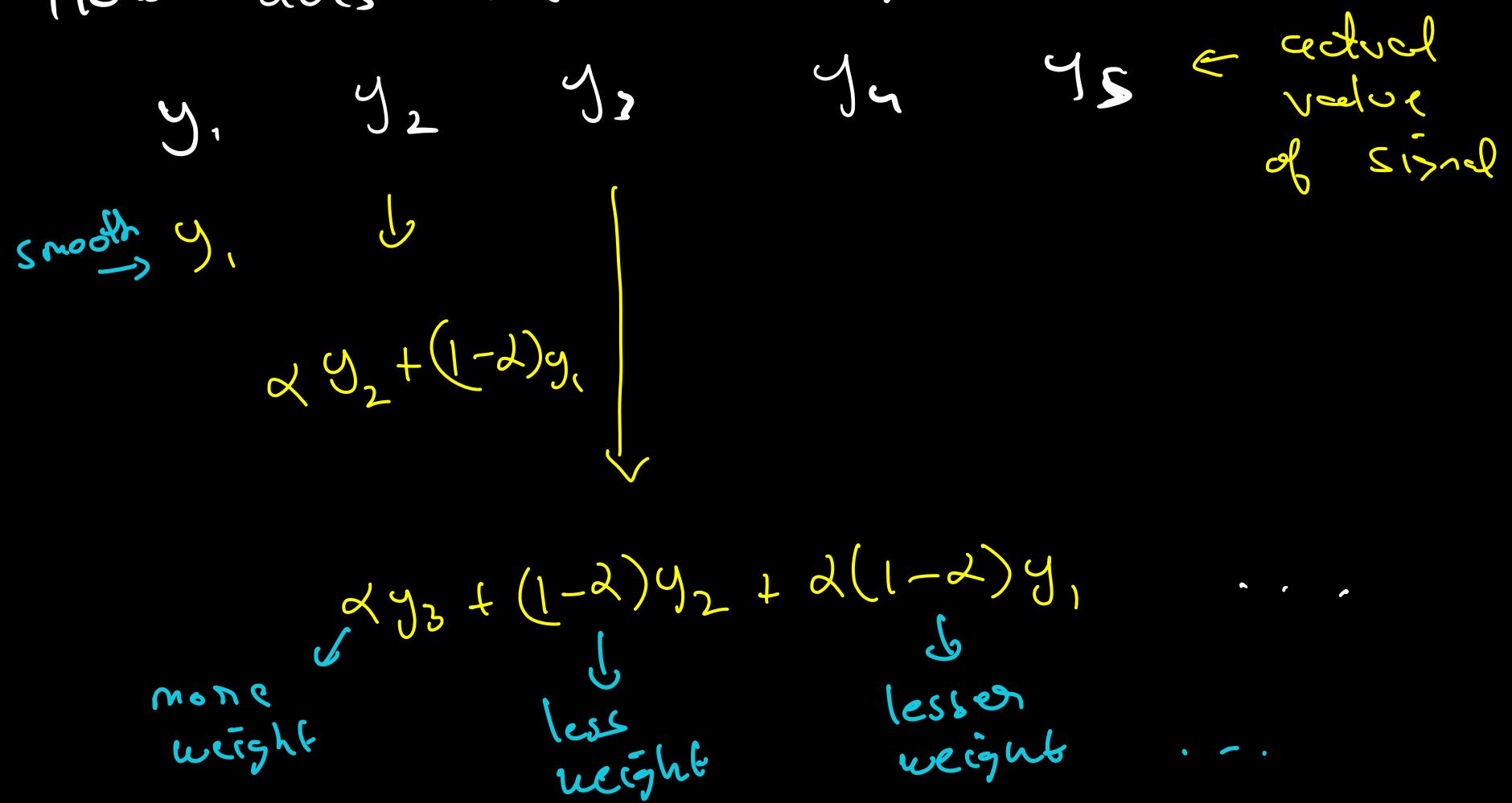
→ Too many oscillations, how to reduce these?

## Simple Exponential smoothing



A.K.A  
Exponential  
moving average

How does this work?



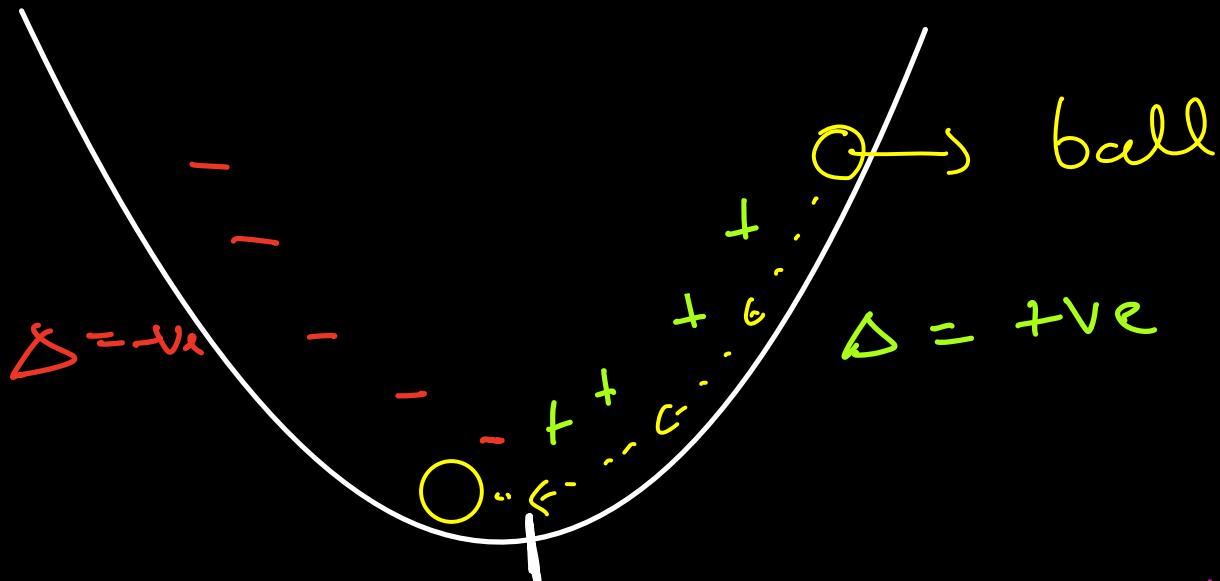
Current value is influenced by past values!

Where I am now + where I have been before?

Recursive Eqn

$$\begin{aligned}\hat{y}_t &= \beta y_t + (1-\beta) \hat{y}_{t-1} \\ &= \beta y_t + (1-\beta) (y_{t-1} + \beta(\hat{y}_{t-2})) \\ &= \beta y_t + (1-\beta)(y_{t-1}) + \beta(1-\beta)(\hat{y}_{t-2}) \\ &\quad \dots \\ &\quad \curvearrowleft \\ &= \beta(1-\beta)y_{t-2} + \beta(1-\beta)^2 \hat{y}_{t-3} \\ &\quad \dots\end{aligned}$$

## Momentum based GW

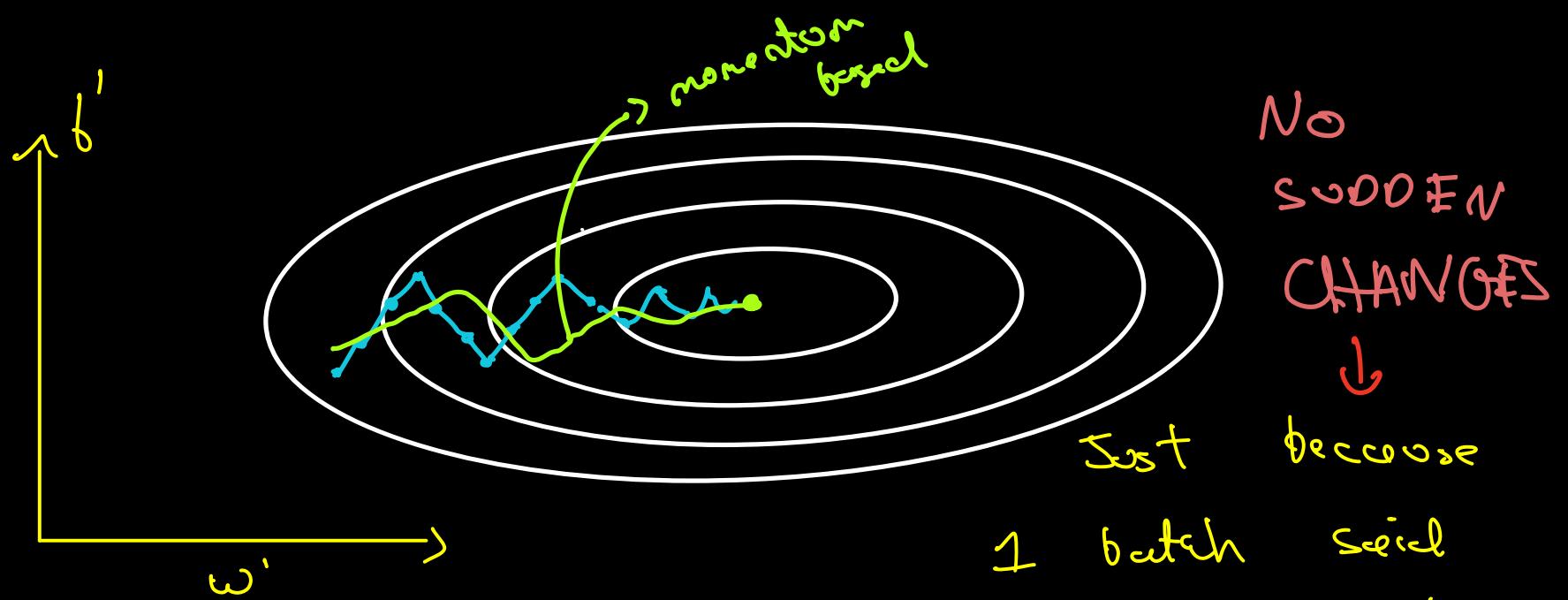


Q: Does the ball change direction immediately? (Physics)

→ No, due to momentum

It was going in the same direction for

many steps. Q: Should it change direction suddenly  
think in terms of mini-batch GD



How to accomplish  
this ?

Normal GD:

$$\begin{aligned}\omega^{t+1} &= \omega^t - \alpha \frac{\partial L}{\partial \omega} \\ &= \omega^t - \lambda d\omega\end{aligned}$$

$$b^{t+1} = b^t - \lambda db$$

Momentum based GD

$V_{d\omega^k} = \beta \overbrace{V_{d\omega^{k-1}}}^{\text{velocity}} + (1-\beta) d\omega^k$   $\curvearrowright$  accelerated

$$V_{db^k} = \beta V_{db^{k-1}} + (1-\beta) db^k$$

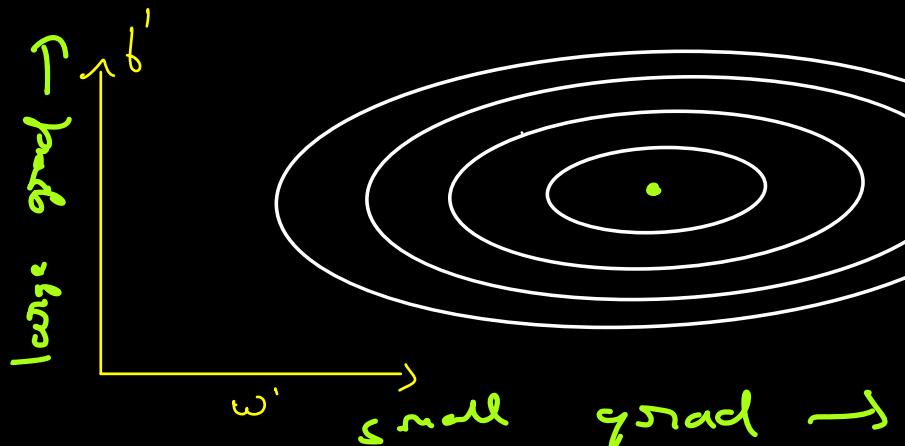
$$\omega^{k+1} = \omega^k - \alpha V_d \omega^k$$

$$b^{k+1} = b^k - \alpha V_d b^k$$

→ code

## RMS prop

there is still one more issue.



$$d\theta \gg d\omega$$

But learning rate is same !

We want to take smaller steps in  
→ direction

$$S_{d\omega^k} = \beta S_{d\omega^{k-1}} + (1-\beta) (d\omega^k)^2$$

$$S_{d\theta^k} = \beta S_{d\theta^{k-1}} + (1-\beta) (d\theta^k)^2$$

still,

$$S_{d\theta} \ggg S_{d\omega}$$

Update  $\rightarrow$

$$\omega^{k+1} = \omega^k - \alpha \frac{\nabla \omega^k}{\sqrt{\nabla f(\omega^k) + \xi}} \rightarrow \text{To avoid div by 0}$$

$$b^{k+1} = b^k - \alpha \frac{\nabla f(\omega^k)}{\sqrt{\nabla f(\omega^k) + \xi}}$$

Now,  
Step size will  
be more limited!!

$\nabla f(\omega^k)$  is divided by  
a large number, so  
it becomes small or

## Adaptive Momentum Estimation (ADAM)

Combination of both !

$$V_{d\omega^k} = \beta_1 V_{d\omega^{k-1}} + (1-\beta_1) d\omega^k$$

$$V_{db^k} = \beta_1 V_{db^{k-1}} + (1-\beta_1) db^k$$

$$S_{d\omega^k} = \beta_2 S_{d\omega^{k-1}} + (1-\beta_2) (d\omega^k)^2$$

$$S_{db^k} = \beta_2 S_{db^{k-1}} + (1-\beta_2) (db^k)^2$$

Update :

$$w^{k+1} = w^k - \alpha \frac{V_{dw^k}}{\sqrt{S_{dw^k} + \epsilon}}$$

$$b^{k+1} = b^k - \alpha \frac{V_{db^k}}{\sqrt{S_{db^k} + \epsilon}}$$

This is most popular and tensorflow default.

Recommended values :

$\alpha \rightarrow$  tune  $\beta_1 = 0.9, \beta_2 = 0.977, \epsilon = 10^{-8}$   
 $\hookrightarrow$  no need to change.

Final note:

Since we are using exp now only,

$$V^{K+1} = \beta V^K + (1 - \beta) dw$$

↓

what is  $V^0 = ?$

→ 0

because of init with 0, for the  
first few steps, there is under estimation  
of momentum, hence there is a small  
modification → Refer Post Recal !