

Optimizers in Neural Network:

NN \rightarrow trained by \rightarrow gradient descent

\rightarrow local minima.

\rightarrow saddle points

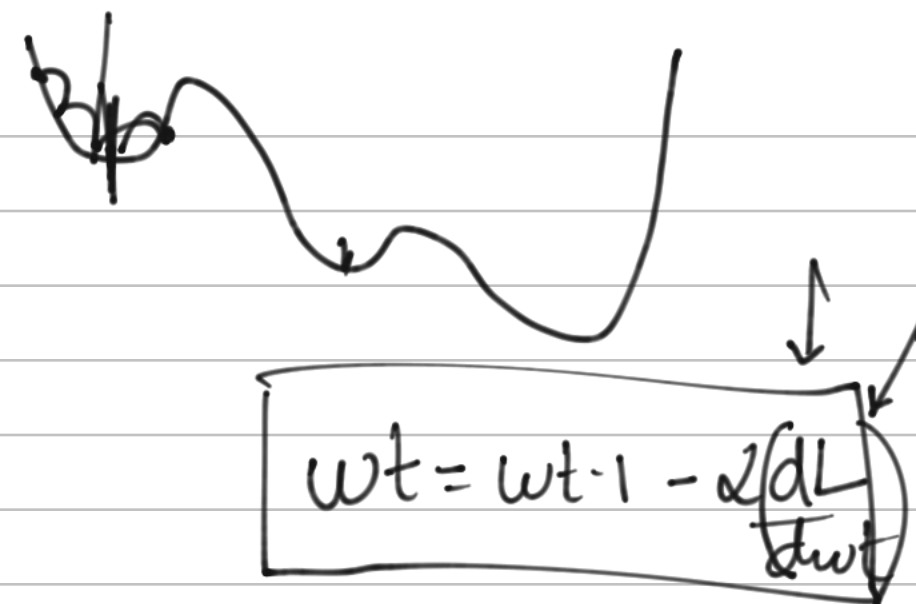
\rightarrow Global min ✓

Optimizers \rightarrow we can atleast try to reach our global minima with more ~~are~~ probability
variants of gd.

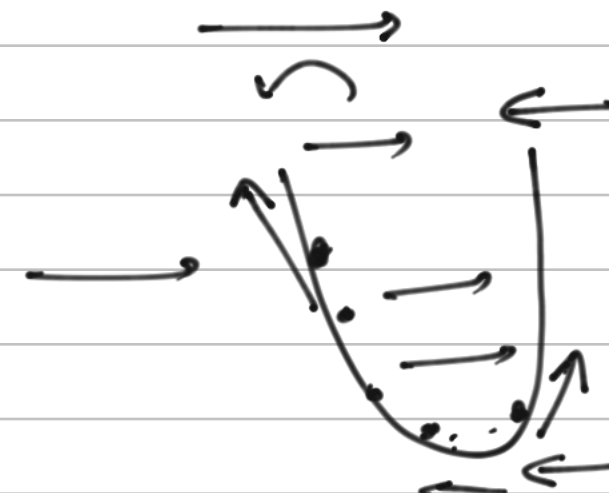
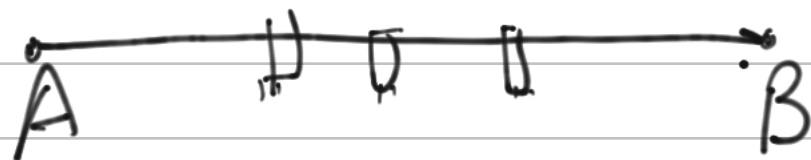
(GD equation) \rightarrow modify \rightarrow

- Momentum based gradient descent
- Adagrad ✓
- Rms Prop ✓
- (Adam) ✓✓

GD



Momentum Based Gradient Descent: optimizer → better GD



Momentum:

GD with momentum

Historical + point grad.

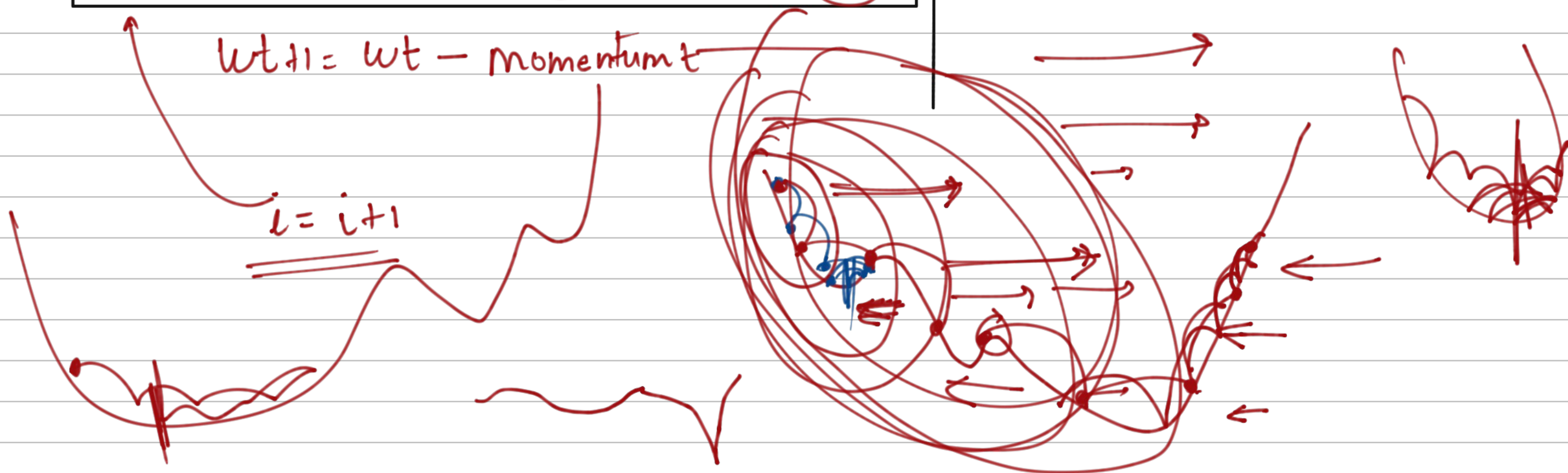
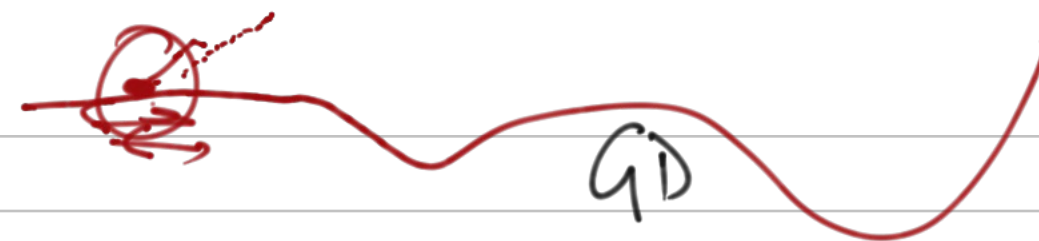
~~Hyperparameter~~

$$\text{momentum}_t = \gamma \text{momentum}_{t-1} + (1-\gamma) \alpha \frac{dL}{dw_t}$$

$$w_{t+1} = w_t - \alpha \frac{dL}{dw_t}$$

$$w_{t+1} = w_t - \text{momentum}_t$$

$$i = i + 1$$

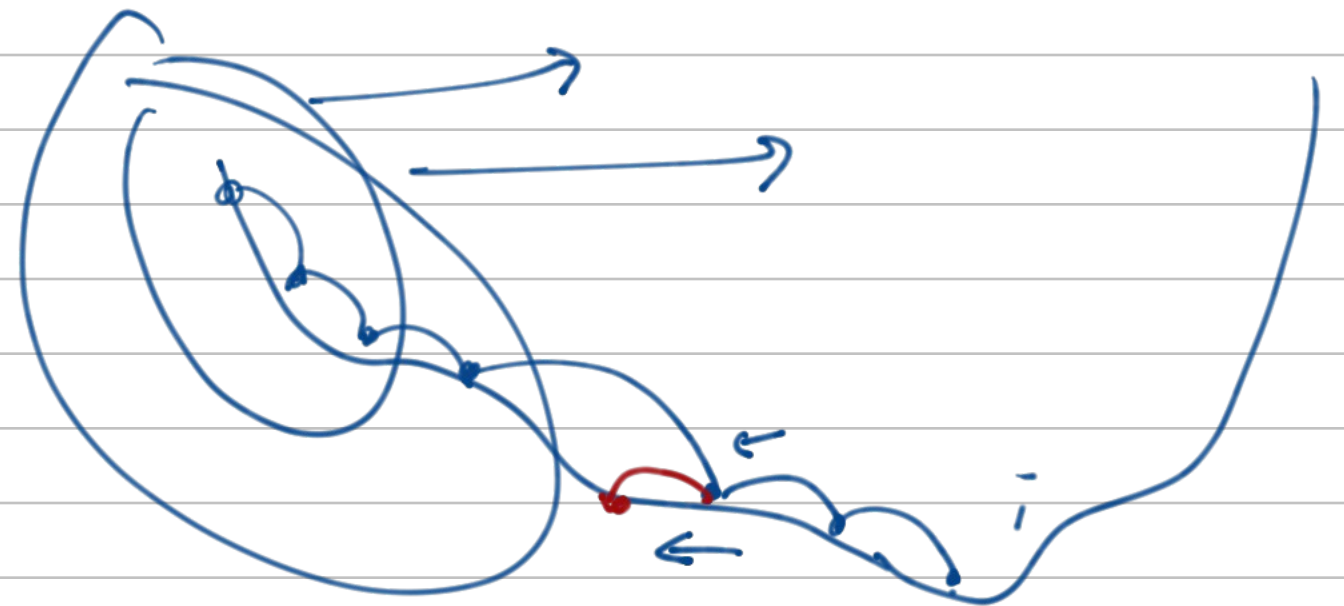


$$\text{momentum}_0 = 0$$

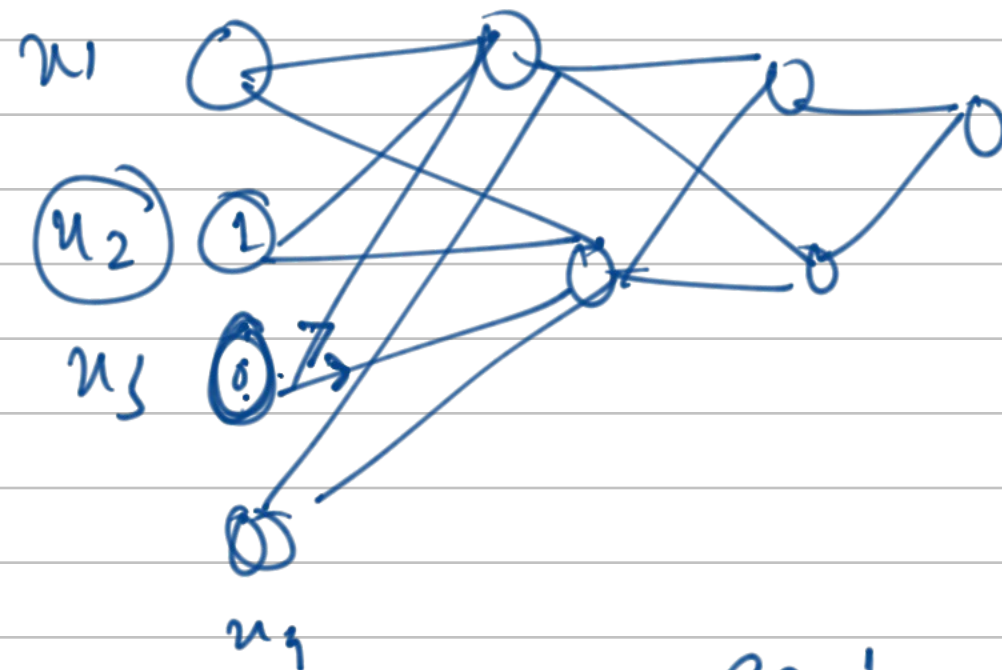
$$\text{momentum}_t = \gamma \text{momentum}_{t-1} + (1-\gamma) \frac{\partial L}{\partial w_t}$$

↓
accumulator
point gradient

$$w_{t+1} = w_t - (\text{momentum}_t)$$



Adagrad \rightarrow Adaptive Gradient \rightarrow optimizer

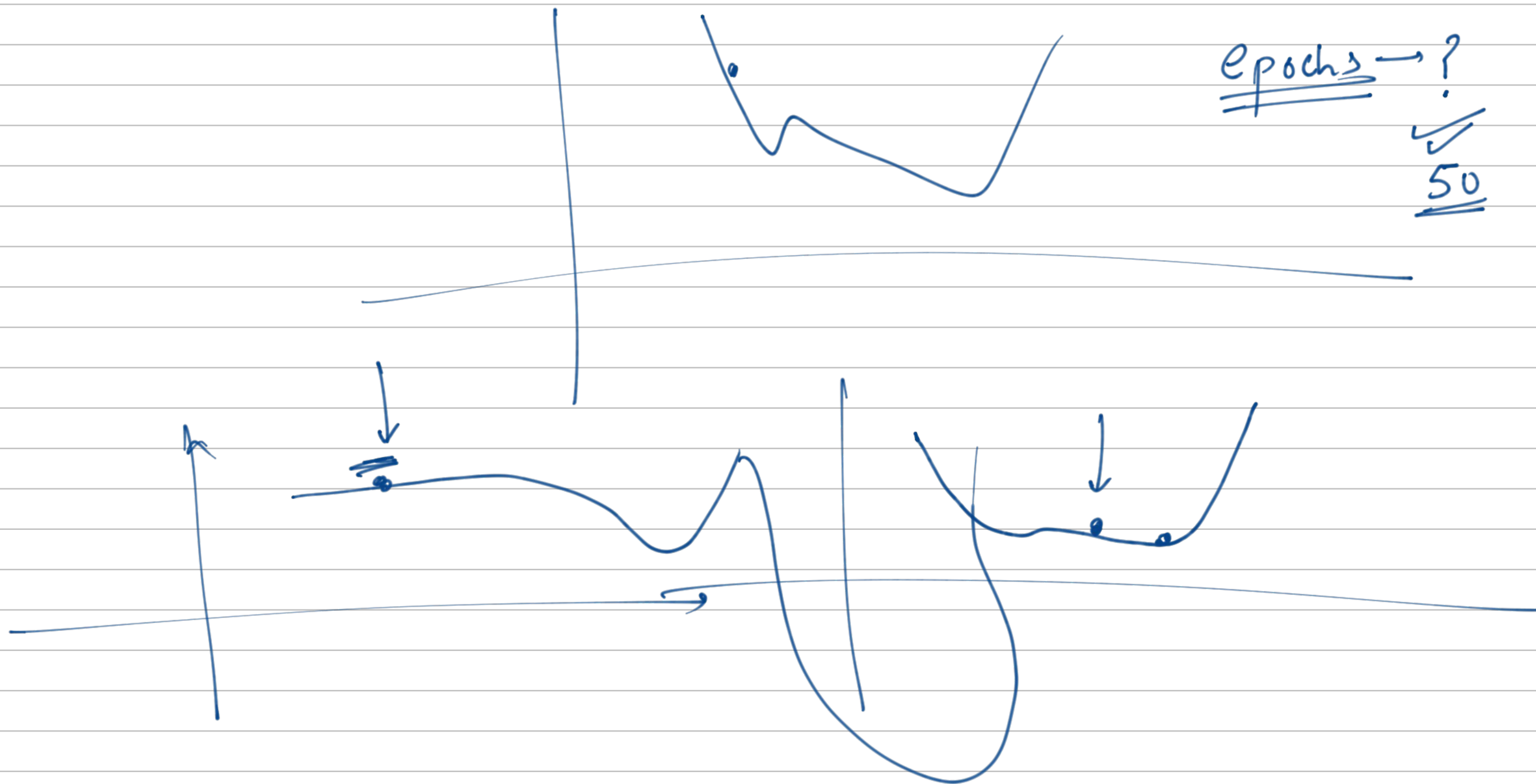
[illegible]

epoch

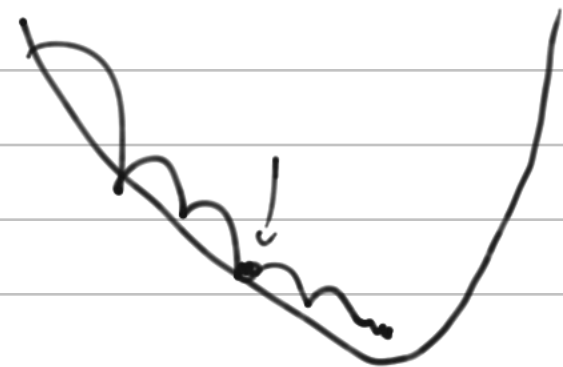
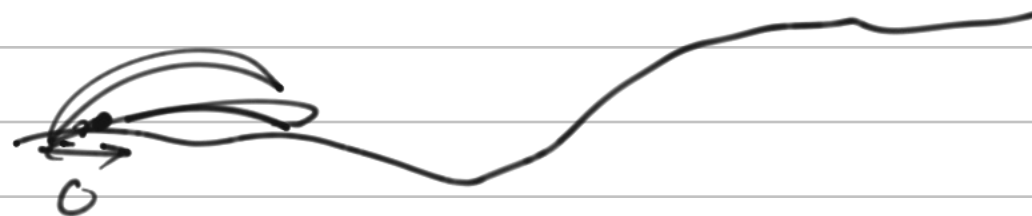
weights associated to u_3 will not have equal opportunity to get trained as much as u_2

All weights are randomly initialized.

epochs → ?
✓
50



Adagrad:



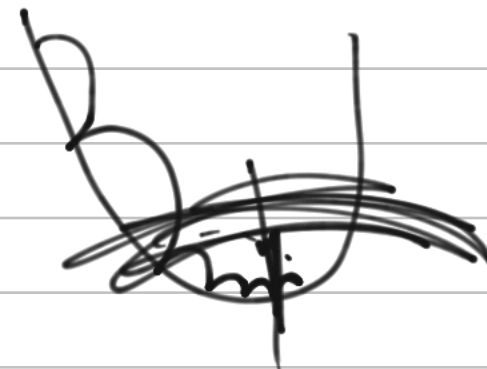
$V_0 = 0$

$$w_{t+1} = w_t - \frac{\alpha}{\sqrt{V_t}} \frac{dL}{dw_t}$$

$$V_t = V_{t-1} + \left(\frac{dL}{dw_t} \right)^2$$

→ Same

people/wts who have
done very big steps



Rms Prop:

$$v_t = \beta v_{t-1} + (1-\beta) \left(\frac{\partial L}{\partial w_t} \right)^2$$

$$w_{t+1} = w_t - \frac{\alpha}{\sqrt{v_{t+1} + \epsilon}} \times \left(\frac{\partial L}{\partial w_t} \right) \rightarrow \text{momentum} - ?$$

Combination of Rms Prop + momentum \rightarrow Adam
Adaptive moments

Adam →

Adaptive moments

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) \times \frac{\partial L}{\partial w_t}$$

momentum

velocity

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) \times \left(\frac{\partial L}{\partial w_t} \right)^2$$

$$w_{t+1} = w_t - \frac{\alpha}{\sqrt{v_t + \epsilon}} \times m_t$$

$$\beta_1 = 0.9 \quad \beta_2 = 0.999$$

