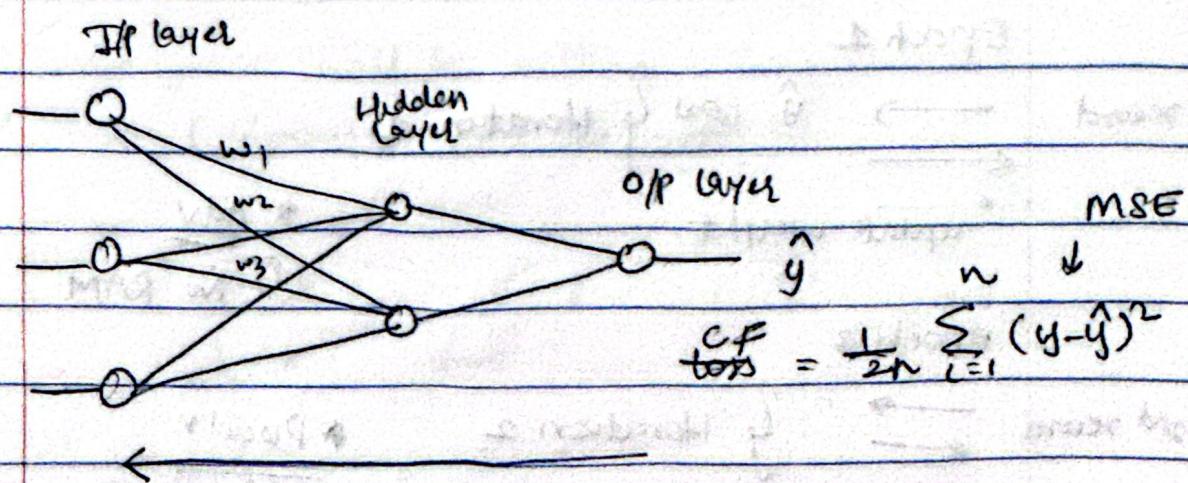
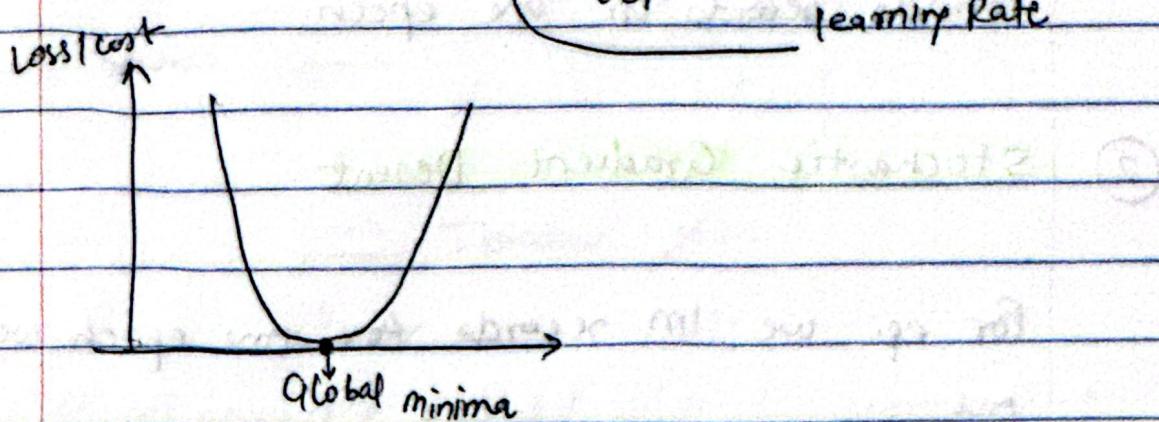


## OPTIMIZERS

### ① Gradient Descent

Weight update formula.

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial L}{\partial w_{\text{old}}} \quad \text{learning rate}$$



To minimize loss funct<sup>n</sup> we use optimizers.

We send all records in gradient descent  
we sent 1 record in epoch is stochastic  
we sent in 1000 batch in mini batch

→ Epoch  $\xrightarrow{100000} \downarrow 1 \text{ Epoch}$

One forward & backward propagation is 1 Epoch.

\* Disadv of gradient

① Resource extensive

(You'll require lots of ram if you have to include 100000 records in one epoch.)

② Stochastic Gradient Descent

For eg. we 1M records for one epoch we pass one

Epoch 1

1 record  $\longrightarrow \uparrow 100 \downarrow \text{Iteration 1}$   
 $\xleftarrow{\text{update weight}}$

\* Adv

① No RAM

Epoch 2

2nd record  $\longrightarrow \downarrow \text{Iteration 2}$  \* Disadv

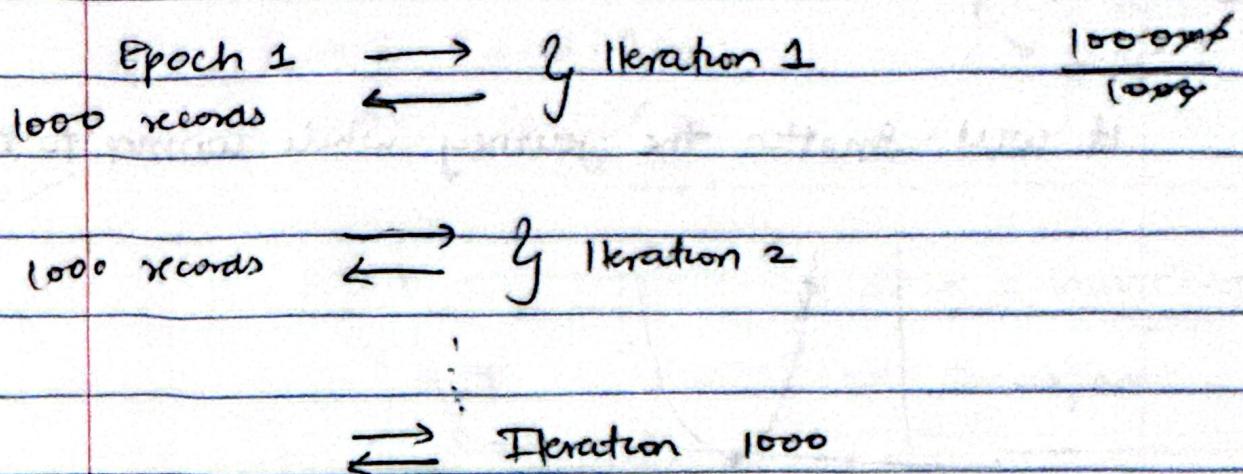
① Convergence will be very slow

Million Iteration

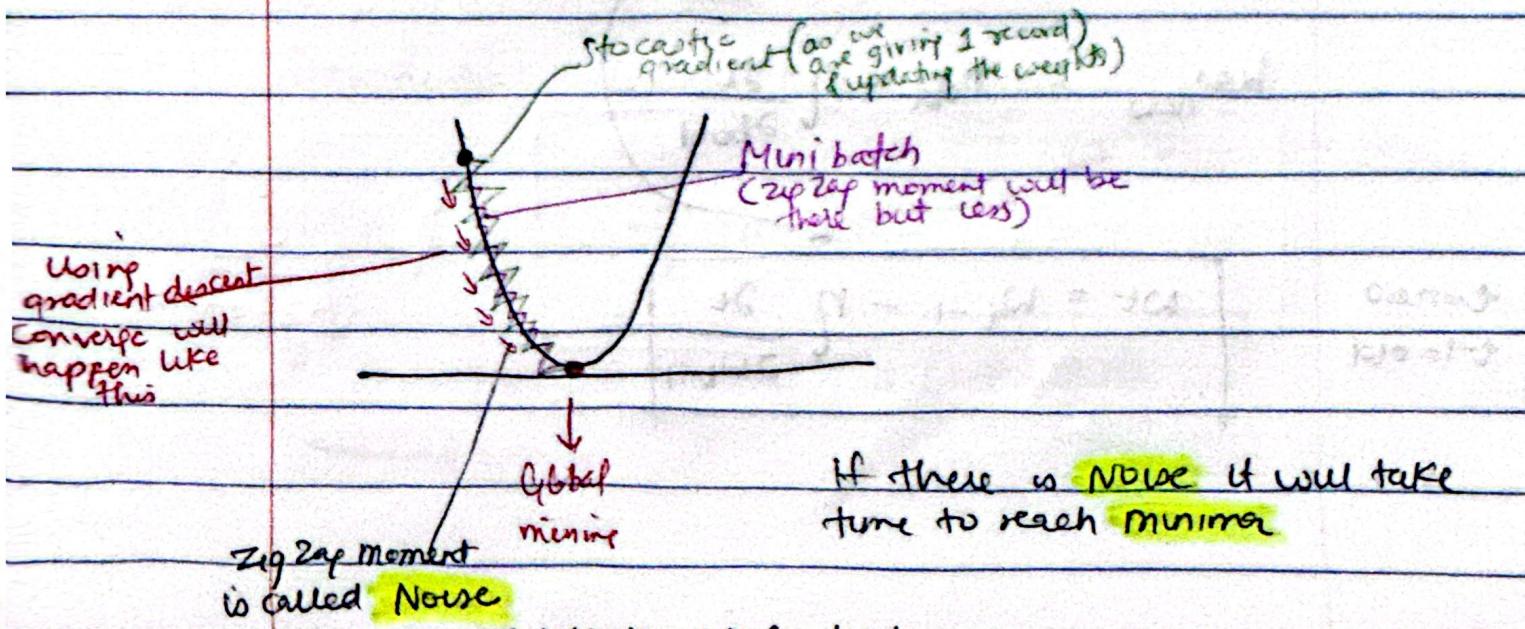
② Time complexity High

### ③ Mini Batch Stochastic Gradient Descent (SGD)

- set a batch size      eg. batch size = 1000



- Not Resource Intensive
- Convergence will be better
- Time Complexity will improve.



\* How do we remove the noise?

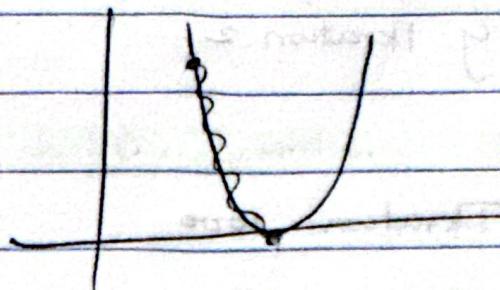
We use concept as **momentum**

(4)

Mini Batch

**SGD with momentum**

It will smoothen the journey while coming to minimum



- Exponential Weighted Avg (moving Average)

Weight update

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial L}{\partial w_{\text{old}}}$$

$$b_{\text{new}} = b_{\text{old}} - \eta \frac{\partial L}{\partial b_{\text{old}}}$$

$t=\text{new}$   
 $t-1=\text{old}$

$$w_t = w_{t-1} - \eta \frac{\partial L}{\partial w_{t-1}}$$

$t = \text{time}$

$$t_1 \quad t_2 \quad t_3 \quad t_4 \dots t_n$$

$$q_1 \quad q_2 \quad q_3 \quad q_4 \dots q_n$$

$\beta \rightarrow \text{Hyperparameter}$

$$v_{t_1} = q_1$$

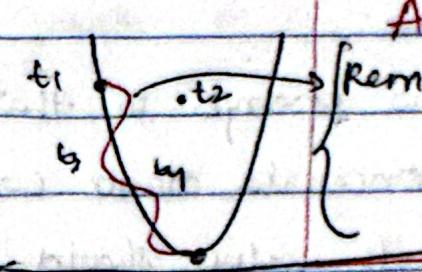
$$\beta = 0 \text{ to } 1$$

$\downarrow$   
0.95

$$v_{t_2} = \beta * v_{t_1} + (1-\beta) * q_2$$

$$= (0.95) * v_{t_1} + (0.05) * q_2$$

weight  
values



Adv  
Removing Noise & Smoothing  
curve.  
Quicker Convergence  
Mini Batch

$$v_{t_3} = \beta * v_{t_2} + (1-\beta) * q_3$$

→ Where do we apply this?

We apply it in  $\frac{\partial L}{\partial w_{t-1}}$

Exponential weighted Avg

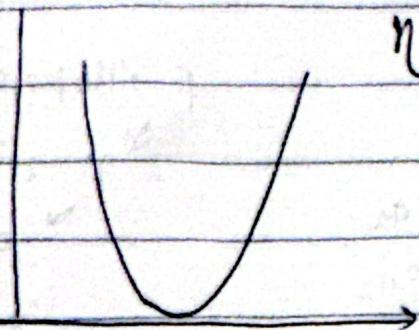
$$v_{t-1} = \beta * v_{t-2} + (1-\beta) * \frac{\partial L}{\partial w_{t-1}}$$

$$w_t = w_{t-1} - \eta (v_{t-1}) \rightarrow$$

$$v_{t-1} = \beta * v_{t-2} + (1-\beta) * \frac{\partial L}{\partial w_{t-1}}$$

## ⑤ Adagrad → Adaptive Gradient Descent

$\eta$  = learning Rate



↳ helps to maintain speed of convergence.

→ Can we somehow change so that if the learning rate is high once it's starts coming close to global minima its value should decrease?

How can we do that

$$w_t = w_{t-1} - \eta \frac{\partial L}{\partial w_{t-1}}$$

Here we replace  $\eta$  with something else as we want to make it adaptive

$$w_t = w_{t-1} - \eta \frac{\partial L}{\partial w_t}$$

eta

$\eta' = \eta / \sqrt{t + \epsilon} \rightsquigarrow$

Small No to avoid dw by 0

$$\alpha_t = \sum_{i=1}^t \left( \frac{\partial L}{\partial w_t} \right)^2$$

↑  
current timestamp

this value will increase so we have  $\Sigma$  (summation)

& thus  $\alpha_t$  is in denominator for  $\eta' = \frac{n}{\sqrt{\alpha_t + \epsilon}}$

so the value of  $\eta'$  will decrease

⇒ By this we are bringing Adaptiveness in the learning rate. So by this the learning rate won't be fixed it will be decreasing as we move towards global minima.

### Disadv

If the value of  $\eta'$  increases  $\alpha_t$  increases too much there will be negligible change of  $\eta'$  so for that we have next optimizer.

We have to make sure  $\alpha_t$  doesn't reach huge value value.

## ⑥ Adadelta & RMS Prop.

$$\eta' = \frac{n}{\sqrt{sdw + \epsilon}}$$

we apply exponential weighted avg.

$$sdw_t = \beta sdw_{t-1} + (1-\beta) \left( \frac{\partial L}{\partial w_{t-1}} \right)^2$$

$$\text{if } \beta = 0.95$$

$$sdw_t = (0.95) sdw_{t-1} + (0.05) \left( \frac{\partial L}{\partial w_t} \right)^2$$

this will control & not make eqn bump up or value to go up.

as a result

$$\eta' = \frac{n}{\sqrt{sdw + \epsilon}}$$

there will be a close decrease in the value.

## Smoothing, learning Rate becomes Adaptive

Smoothng op<sup>n</sup> will have  $V_{dw}$  as in SGD with momentum but here in Adadelta smoothing wont happen,

there is no  $V_{dw}$  here.

Now this  $V_{dw}$  &  $s_{dw}$  should be combined together.

### ⑦ Adam Optimizers (Best Optimizer)

(momentum + RMS Prop (Adaptive Learning Rate))

$$V_{dw} | V_{db}$$

$\downarrow$   
bias

$$s_{dw} | s_{db}$$

$\downarrow$   
bias

Now my weight update op<sup>n</sup> will be,

$$w_t = w_{t-1} - \eta' V_{dw}$$

for Bias update

$$b_t = b_{t-1} - \eta' V_{db}$$

$$\eta' = \frac{\eta}{\sqrt{s_{dw} + \epsilon}}$$

$$V_{db_t} = \beta * V_{db_{t-1}} + \frac{(1-\beta) \partial L}{\partial b_{t-1}}$$

$$V_{dw_t} = \beta * V_{dw_{t-1}} + \frac{(1-\beta) \partial L}{\partial w_{t-1}}$$