# Day 3 - Deep Learning.

① Optimizers

    i) Gradient Descent

    2) SGD (Stochastic Gradient Descent)

    ③ Mini Batch SGD

    ④ SGD with Momentum

    ⑤ Adagrad

    ⑥ RMSPROP

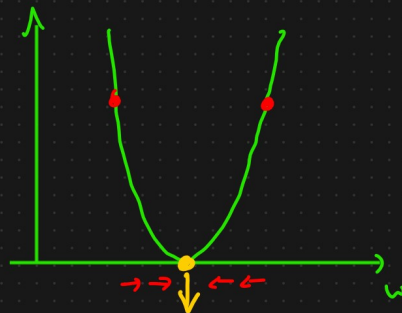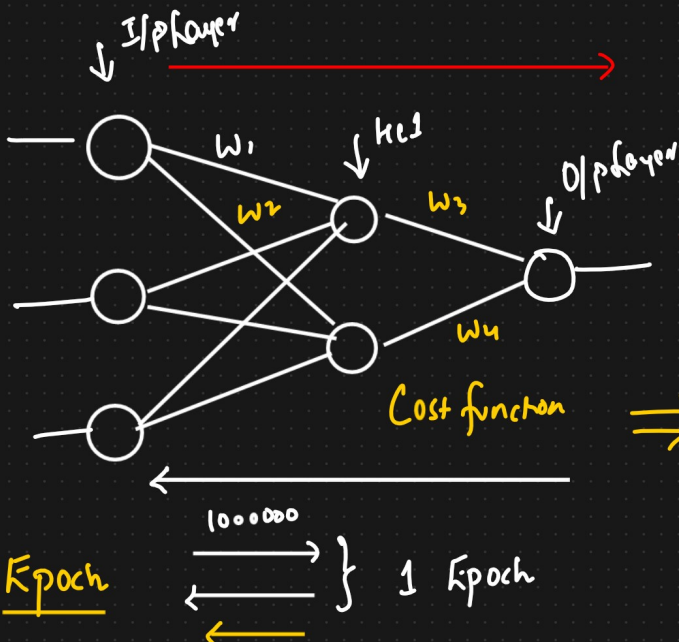    ⑦ Adam optimizer

Batch, Epochs, Iterations

⇓

ANN

---

① GRADIENT DESCENT → Optimizer

Weight Updation Formula → Learning Rate

$$W_{new} = W_{old} - \eta \frac{\partial h}{\partial W_{old}}$$

Loss/Cost

Global Minima.

I/p layer

$W_1$    Hid    O/p layer

$W_2$   $W_3$   $\hat{y}$

$W_4$

Cost function

MSE    Optimizers

$$\Rightarrow \frac{1}{2n} \sum_{i=1}^{n} (y - \hat{y})^2 \}$$

1000000

Epoch  ⇄ } 1 Epoch

{ 1000000 }

## Disadvantage

① Resource Extensive {huge RAM}

② Stochastic Gradient Descent

Epoch 1      → 100 epochs

1000000

  ① RAM ↓↓

1 record ⟶ $\hat{y}$ } → Iteration 1
⟵
⟵ loss
Update weights

  Disadvantage

ⓕ Convergence will } be very slow

2 record ⟶ } Iteration 2
⟵

ⓕ Time complexity will also be high

} Million Iteration

$\dfrac{1000\cancel{\phi\phi}}{1\cancel{\phi\phi}}$

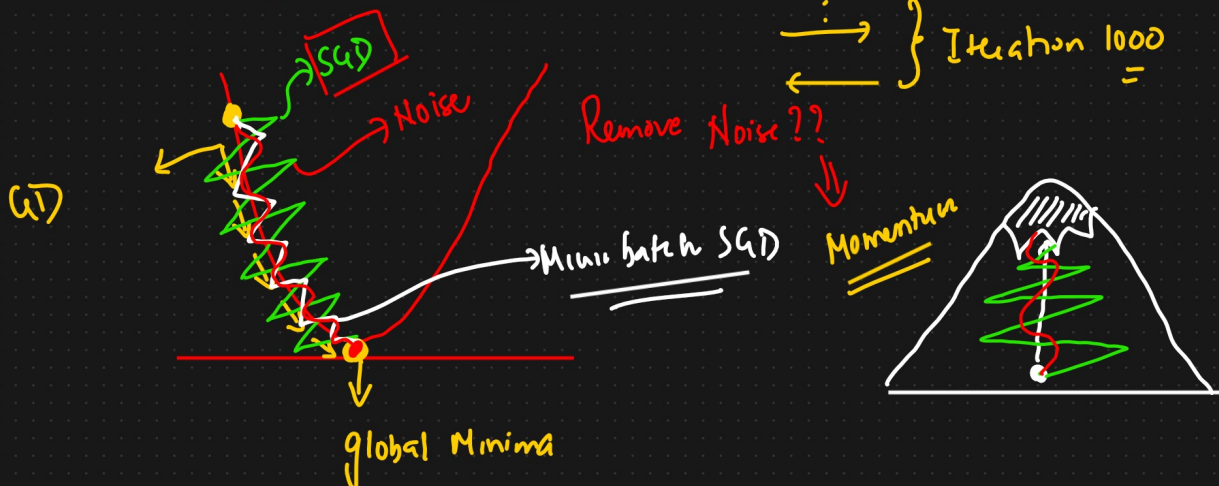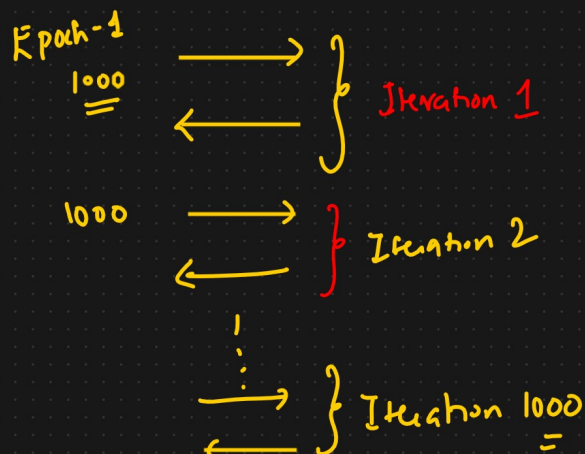③ Mini batch SGD

1000000     batch-size = 1000

Ⓐ Resource Intensive

② Convergence will be better

③ Time Complexity will Improve

Epoch-1   ⟶ } Iteration 1
1000 ⟵

1000 ⟶ } Iteration 2
⟵

⟶ } Iteration 1000
⟵

SGD
→ Noise
GD

Remove Noise ??

→ Mini batch SGD   Momentum

global Minima

# ④ SGD With Momentum

{ Exponential Weighted Average }

$$W_{new} = W_{old} - \eta \frac{\partial h}{\partial w_{old}}$$

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

⇓
Time Series
⇓
ARIMA, ARMA,

$$\boxed{W_t = W_{t-1} - \eta \frac{\partial h}{\partial w_{t-1}}}$$

## Exponential Weighted Average

{ Forecasting }

$$t_1 \quad t_2 \quad t_3 \quad t_4 \quad \cdots \quad t_n$$
$$a_1 \quad a_2 \quad a_3 \quad a_4 \quad \cdots \quad a_n$$

$\beta \Rightarrow$ Hyper parameter

$\beta = 0 \text{ to } 1$
⇓
0.95

$$V_{t_1} = a_1$$

$$V_{t_2} = \beta * V_{t_1} + (1-\beta) * a_2$$

$$= (0.95) * V_{t_1} + (0.05) * a_2$$

→ Removing noise
smoothing
the curve

$$V_{t_3} = \beta \times V_{t_2} + (1-\beta) * a_3$$

## Exponential Weighted Avg ←

$$W_t = W_{t-1} - \eta \, V_{dw}$$

⇑

$$\boxed{V_{dw_t} = \beta \times V_{dw_{t-1}} + (1-\beta) * \frac{dL}{\partial W_{t-1}}}$$

{ Reduce the noise  
Minibatch  
Quicker Convergence

## Recap

① Gradient Descent

② SGD

③ Mini batch SGD

④ SGD with Momentum

⑤ Adagrad → Adaptive Gradient Descent

{ fixed } ⇒ optimiza  
↑  
$\eta$ = Learning Rate



$\eta$ = fixed ⇒ adaptive ⇒ Learning Rate ⇒ Decreasing → Global Minima

Global Minima ↓

$$W_t = W_{t-1} - \boxed{\eta} \frac{dL}{\partial w_{t-1}}$$

↓↓↓

$$W_t = W_{t-1} - \eta' \frac{\partial L}{\partial w_{t-1}} \checkmark$$

$$W_t \approx W_{t-1}$$

↓↓↓ $\eta' = \eta \swarrow 0.01$ ϵ

$$\Rightarrow \sqrt{\alpha_t + \epsilon} \rightarrow$$ Small number

$\eta \searrow$ ↓

↓ Decreasing ↓↓↓

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

Huge number

| $t=1$ | $t=2$ | $t=3$ |
|---|---|---|
| $\eta = 0.01$ | $\eta = 0.005$ | $\eta = 0.002$ |

⑥ Adadelta And RMSProp

Exponential Weighted Average  
↓

$$\eta' = \frac{\eta}{\sqrt{Sdw + \epsilon}} \Bigg\}$$

$Sdw_{t-1} = 0$

$$Sdw_t = \beta\, Sd\overset{v}{w}_{t-1} + (1-\beta)\left(\frac{\partial h}{\partial w}\right)^2_{t-1}$$

$$\beta = 0.95$$

$$\boxed{Sdw_t} = (0.95)\, Sdw_{t-1} + (0.05)\left(\frac{\partial h}{\partial w_{t-1}}\right)^2$$

---

(★) **Adam Optimizer** (Best Optimizer)

Momentum + RMSPROP (Adaptive Learning Rate) $\Bigg\{$ ① Smoothening ② Learning Rate Adaptive $\Bigg\}$

$Vdw=0 \quad Vdb=0 \quad Sdw=0 \qquad Sdb=0$

$$\boxed{\begin{array}{l} W_t = W_{t-1} - \eta'\, Vdw \\[6pt] b_t = W_{bt-1} - \eta'\, Vdb \end{array}}$$

$$\eta' = \frac{\eta}{\sqrt{Sd_w + \epsilon}}$$

$$\boxed{Vdw_t = \beta \times Vdw_{t-1} + (1-\beta)\frac{\partial h}{\partial w_{t-1}}}$$

$$\boxed{Vdb_t = \beta \times Vdb_{t-1} + (1-\beta)\frac{\partial h}{\partial b_{t-1}}}$$