# Deep Learning — Optimizers

## Optimizers
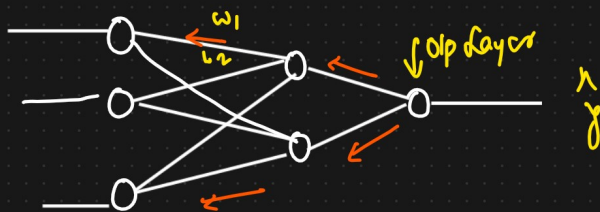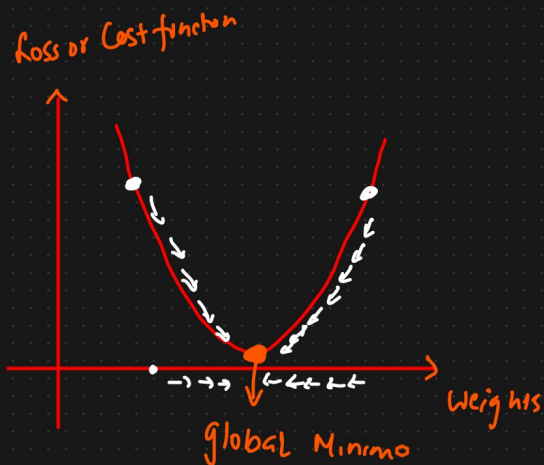
1. Gradient Descent } ✓
2. SGD } ✓
3. Mini batch SGD } ✓
4. SGD With Momentum
5. Adagrad And RMSPROP
6. Adam Optimizers

## Ⓚ GRADIENT DESCENT

### Weight updation Formula

$$W_{new} = W_{old} - \eta \left( \frac{\partial w}{\partial w_{old}} \right)$$

→ Learning Rate

Loss or Cost function



Weights

Global Minimo

$w_1$
$L_2$

↘ O/p layer

$\hat{y}$

Optimizers

Loss fn ↓  ⟹ Backpropogation
or Cost fn ↓

$n = 10000$

### MSE

$$loss\,fn = (y - \hat{y})^2$$  → 1 Records

$$Cost\,fn = \frac{1}{n} \sum_{i=1}^{n} (y - \hat{y})^2$$ → n Records

Epochs, Ituation

$1$ Epoch
$10000$ }  →  $\hat{y}$ ⟹ Cost function ↓

weight will get update

Disadvantage

$n = 1000000$    64
8gh

① Resource Intensive    { Huge Ram, High class GPU }    10k * 100

$n = 1000000$    $n = 10000$

② Stochastic GRADIENT DESCENT    Epoch 1

1 Epoch = 10k Iteration

2 Epoch = 10k

Advantage

① { Solve Resource Issue }    1 record
}→ Iteration 1    $loss = (y - \hat{y})^2$

Disadvantage    2 record }
3rd record }    100 Epochs
10K * 100

① Time Complexity

② Convergence will be very slow.


Noise

$n = 10000$    100 Epochs    Epoch 1

SGD

1 record
loss } Iteration 1

1 Epoch = 10k Iteration

100 Epochs = ?    2nd record
loss } Iteration 2

100 * 10K ← Algebra    ⋮

③ Mini batch SGD    Huge Resources

batch size    SGD ⇌ Slow Convergence

Cost function
↓    Advantage

$n = 100000$    batch-size = 1000    ⊛ Convergence speed will increase

Epoch 1    $\frac{1000}{100} = 100$ Iteration

100 records    ⊛ Noise will be less

$\Downarrow$

Mini batch SGD $\rightarrow$ Noise is there

Convergence become **faster**

① SGD with Momentum

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

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

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

{Exponential Weighted Average}

$\Downarrow$

ARIMA, ARMA

$\Downarrow$

Time Series

Exponential Weighted Average    {Smoothening}

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

$$\boxed{a_1 \quad a_2} \quad a_3 \quad a_4 \quad - - - a_n$$

$\beta :$ Hyperparameter

$$V_{t_1} = a_1$$

$\beta :$ 0 to 1

$$\boxed{\beta : 0.95} \leftarrow$$

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

$$= \beta * a_1 + (1-\beta) * a_2$$

$$= 0.95 * a_1 + 0.05 * a_2$$

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

$$= \beta \left[ 0.95 \, a_1 + 0.05 * a_2 \right] + (1-\beta) * a_3$$

$$= 0.95 \left( 0.95 \, a_1 + 0.05 * a_2 \right) + (0.05) * a_3.$$

## Advantage

Ⓐ Reduces the noise  
Ⓑ Quick Convergence } .

## Recap

① Gradient Descent [ Rich ] → 1 Epoch = 1 Iteration

② ↑ SGD → 1 Epoch = n Iteration

③ ↑ Mini batch SGD { Noise } → 1 Epoch = datsize/batch-size

④ SGD with Momentum

⑤ Adagrad : Adaptive GRADIENT DESCENT



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

$\boxed{\eta = fixed}$

dynamic value

$$W_t = W_{t-1} - \eta' \sqrt{\frac{\partial h}{\partial W_{t-1}}}$$

Initial

$$\eta' = \frac{\acute{\eta} \longrightarrow \text{learning Rate}}{\sqrt{\alpha_t + \epsilon}} \uparrow\uparrow\uparrow \quad \epsilon \rightarrow \text{small value}$$

$$\boxed{W_t \approx W_{t-1}} \Leftarrow \qquad \alpha_t = \sum_{i=1}^{t} \left( \frac{\partial h}{\partial W_t} \right)^2$$

$t=1 \qquad t=2 \qquad t=3 \qquad \cdots \cdots$ <u>Convergence</u>

$\eta = 0.01 \qquad \eta = 0.005 \quad \eta = 0.003$

<u>Exponential Weight Average</u>

⑥ <u>Adadelta And RMSPROP</u>

$t=1 \qquad Sd_{w_t} = 0 \qquad \{ \text{Dynamic } \ell R + \text{Smoothening} \}$

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

$t=2 \qquad Sd_{w_t} = \beta * Sd_{w_{t-1}} + (1-\beta) \left( \frac{\partial h}{\partial w_{t-1}} \right)^2$

$t=3$
$\beta = 0.95$

$$W_t = W_{t-1} + \eta' \frac{\partial h}{\partial w_{t-1}}$$

⑦ <u>Adam optimizer</u> [You should Akam]

<u>SGD with Momentum + RMSPROP</u> [Dynamic $\ell R$ + Smoothening]

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

$$\eta' = \frac{\eta}{\sqrt{S_{dw} + \epsilon}} \quad \Leftarrow$$

$V_{dw} = 0$

<u>Smoothening</u>

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

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