

Day-3.

* Optimizers.

- ① Gradient Descent
- ② SGD (Stochastic Gradient D).
- ③ Mini Batch SGD
- ④ SGD with momentum
- ⑤ Adagrad.
- ⑥ RMSPROP
- ⑦ Adam optimizers.

* Batch
 * Epochs.
 * Iterations.

} ANN.

① Gradient Descent ← optimizers.

* Disadvantage of Gradient Descent

- ① Resource Extensive (Huge Ram)

* Stochastic Gradient Descent

1 record. $\xrightarrow{\quad} \hat{y} \xleftarrow{\quad}$ } → Iteration 1.

2 record. $\xrightarrow{\quad} \hat{y} \xleftarrow{\quad}$ } → Iteration 2

* Disadvantage of Stochastic Gradient Descent.

① Convergence will be very slow.

* Mini Batch (SBD).

Epoch 1 \longleftrightarrow } Iteration 1.
10000

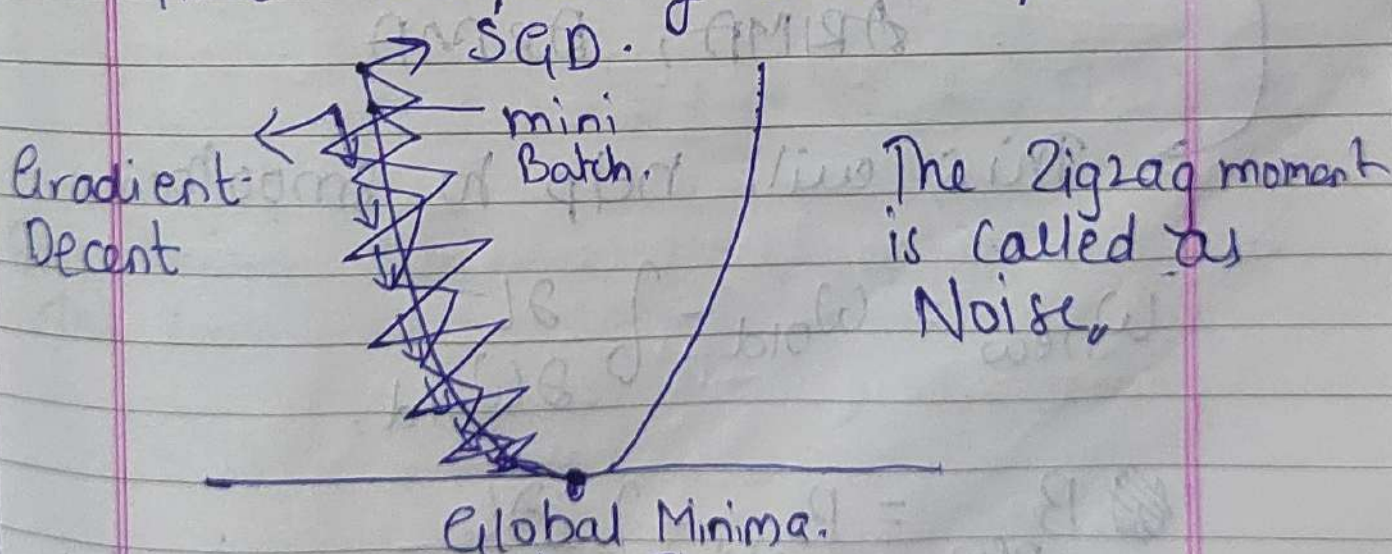
Epoch 2 \longleftrightarrow } Iteration 2

} \rightarrow Iteration 3.

\rightarrow Resource Intensive.

\rightarrow Convergence will be better.

\rightarrow Time complexity will improve.



* In Mini Batch the Zigzag moment will be less.

* In SGD the Zigzag will be high.

* There is a Noise

How do we remove the noise??

→ We use a concept which is called as momentum.

④ SD with momentum

→ What is this momentum will do is that it will smoothen this journey

{ Exponential Weight Average }

↓

Time Series

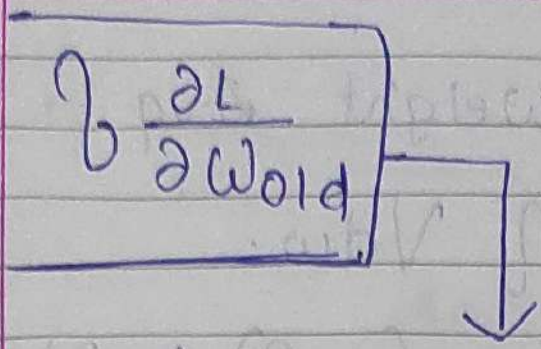
↓

ARIMA, ARMA

→ This will help to smoothen

$$W_{\text{new}} = W_{\text{old}} - \int \frac{\partial L}{\partial W_{\text{old}}}$$

$$B_{\text{new}} = B_{\text{old}} - \int \frac{\partial L}{\partial B_{\text{old}}}$$



w = weights.
 t = Time.

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

* Exponential weights Average.

Data set:

t_1	t_2	t_3	t_4	\dots	t_n
a_1	a_2	a_3	a_4	\dots	a_n

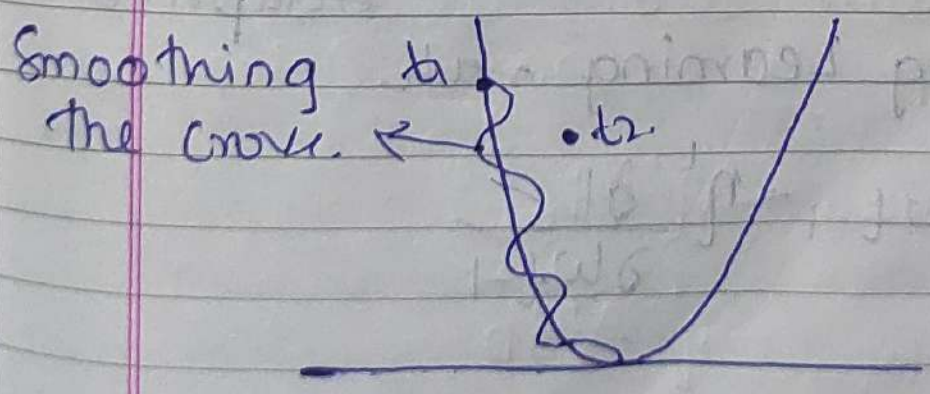
V = Value.

β = Beta. if $\beta = 0.95$

$$V_{t_1} = a_1$$

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

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

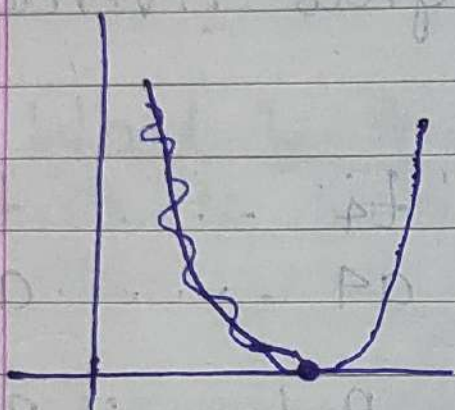


Exponential Weight Avg.

$$w_t = w_{t-1} - \eta \nabla dw$$

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

* The Problem we're solving.



1. Reducing the Noise.

2. Mini Batch.

3. Quick Convergence.

* AdaGrad \rightarrow Adaptive.

Gradient Descent.

$$w_t = w_{t-1} - \boxed{\eta} \frac{dL}{dw_{t-1}} \quad \eta = \text{fixed}$$

\Downarrow

Adaptive

* Changing Learning rate

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

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

$\epsilon = \text{Epsilon}$

↓

→ small number

↓

denominator never
become 0

* Our main aim is to keep decreasing as we reach global minima.

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

$t = \text{Current time-stamp}$

↑↑↑

$$\begin{array}{ccc} t_1 & t_2 & t_3 \\ \eta = 0.01 & 0.05 & 0.02 \end{array}$$

* As we reach near global minima this value will keep on decreasing that why we are bringing Adaphiveness in the learning rate.

* That why the learning rate never be fixed. it will be decreasing as we move towards global minima.

$\Delta t =$ will be a Huge Number.

alpha does not reaches a Huge value
to prevent this.

Page No.

Date: | |

↓
* Adadelta & RMSProp.

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

$$s_{dw} = \beta \cdot$$

$$s_{dw_{t-1}} + (1 - \beta) \left(\frac{\partial L}{\partial w_{t-1}} \right)^2 \quad \text{(EWA)}$$

$$\beta = 0.95$$

$$s_{dw_t} = (0.95) s_{dw_{t-1}} + (0.05) \left(\frac{\partial L}{\partial w_{t-1}} \right)^2$$

* This will always increases with the smaller number because of Beta value.

* Adam

In Adam Optimizer we combine momentum along with the RMSprop it becomes Adam Opt

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

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

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

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

* It is solving the problem of !

① Smoothing.

② Learning rate becomes Adaptive.