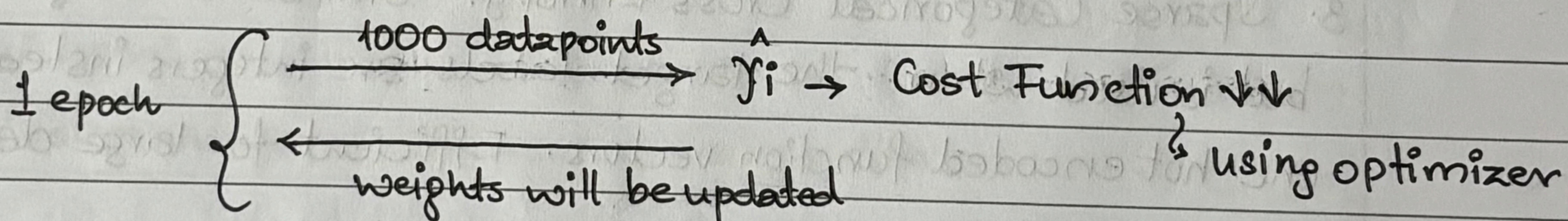


## ⑦ \*Optimizers\*

An optimizer improves the model by adjusting its parameters (weights and biases) to minimize the loss function value.

The optimizers role is to find the best combination of weights and biases that leads to the most accurate predictions.

→ Epochs: One complete pass of the entire training dataset through the learning algorithm



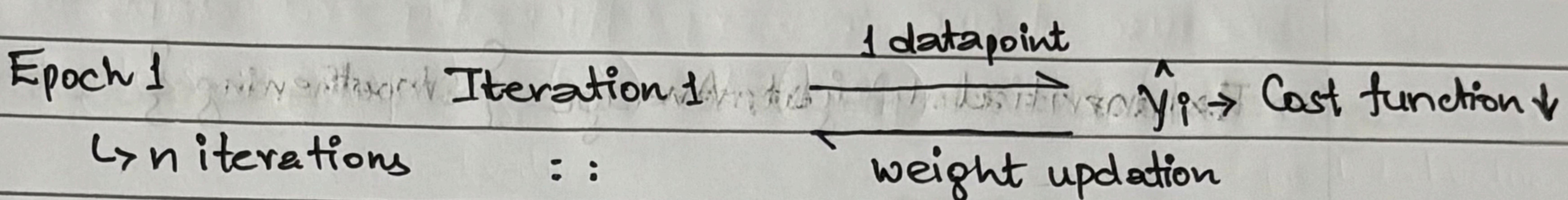
[1 epoch = 1 iteration]

1. Gradient Descent: It works by iteratively adjusting the model parameters in the direction that minimises the loss function.

⇒ Advantages: Convergence will happen

⇒ Disadvantages: High Resource (RAM, GPU) needed

2. Stochastic Gradient Descent : SGD updates the model parameters after each training example, making it more efficient for large datasets compared to traditional Gradient Descent, which uses the entire dataset for each update.



→ Advantages: It solves resource issue

→ Disadvantages:

1. Time Complexity increases, convergence will also take time
2. Noise gets introduced

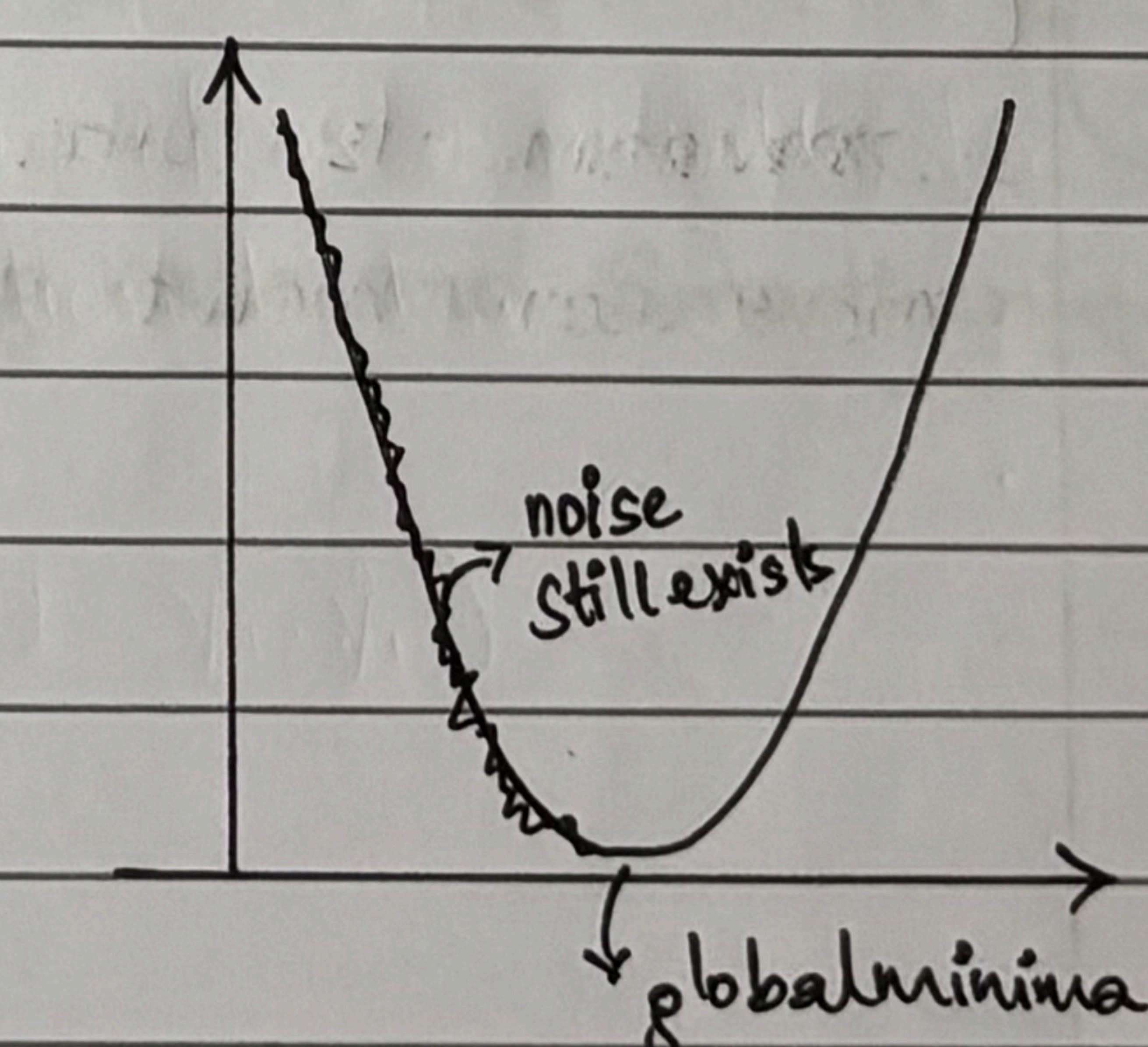
3. Mini Batch Gradient Descent: consists of predetermined number of training examples, smaller than the full dataset. This approach combines the advantages and features of the previous variants.

Datapoints = 100,000      Batch Size = 1000

$$\text{No of iterations} = \frac{\text{Datapoints}}{\text{BatchSize}} = \frac{100,000}{1000} = 100 \text{ iterations}$$

(1 Epoch = No of Iterations) → Noise will be reduced

→ Disadvantages: Noise still exists



#### 4. Stochastic Gradient Descent with Momentum:

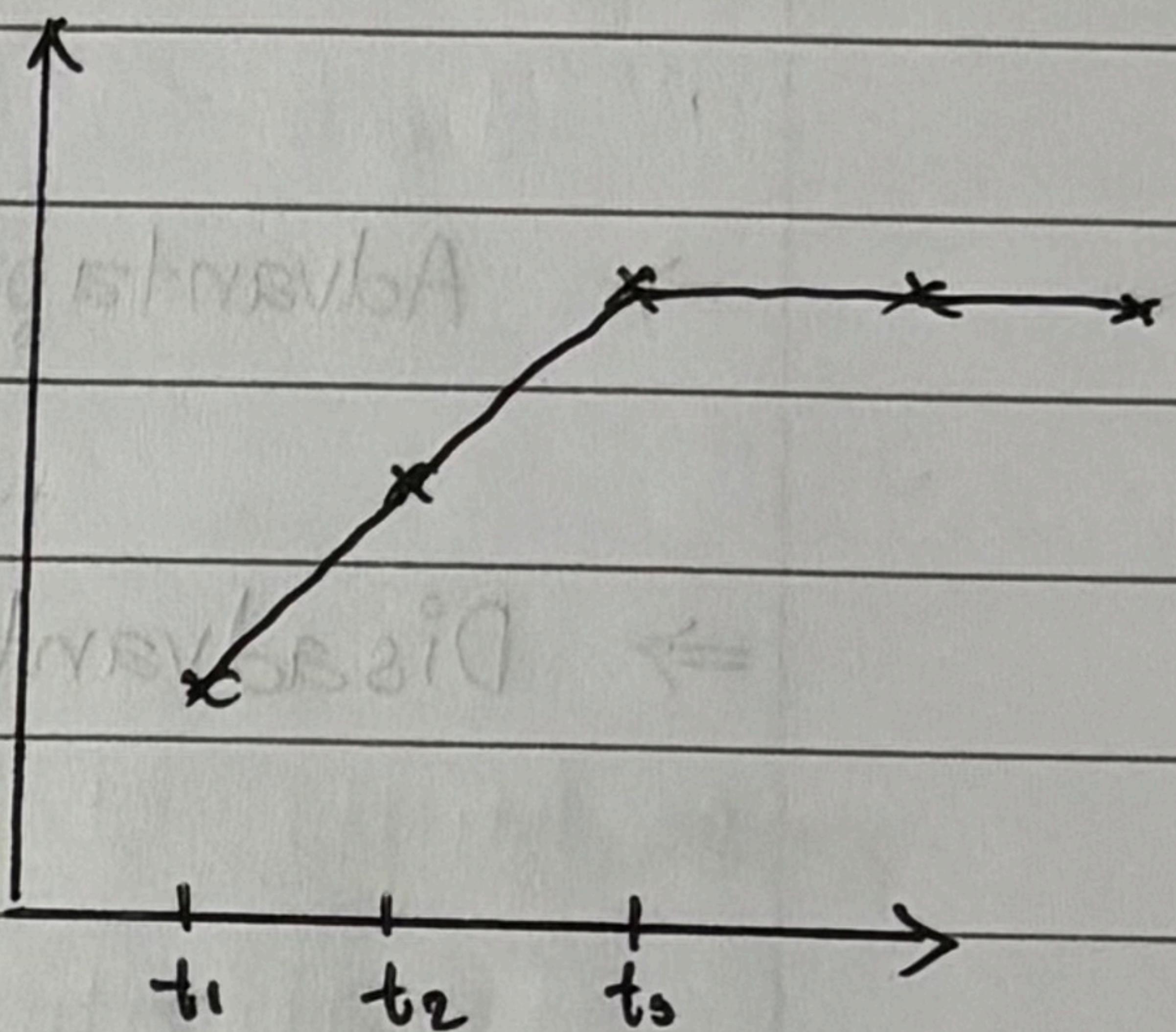
Momentum helps accelerate convergence by smoothing out the noisy gradients of SGD, thus reducing fluctuations and improving speed of convergence.

→ Exponential Weight Average { Smoothening }

Suppose Time =  $t_1, t_2, t_3, t_4, \dots, t_n$

Values =  $a_1, a_2, a_3, a_4, \dots, a_n$

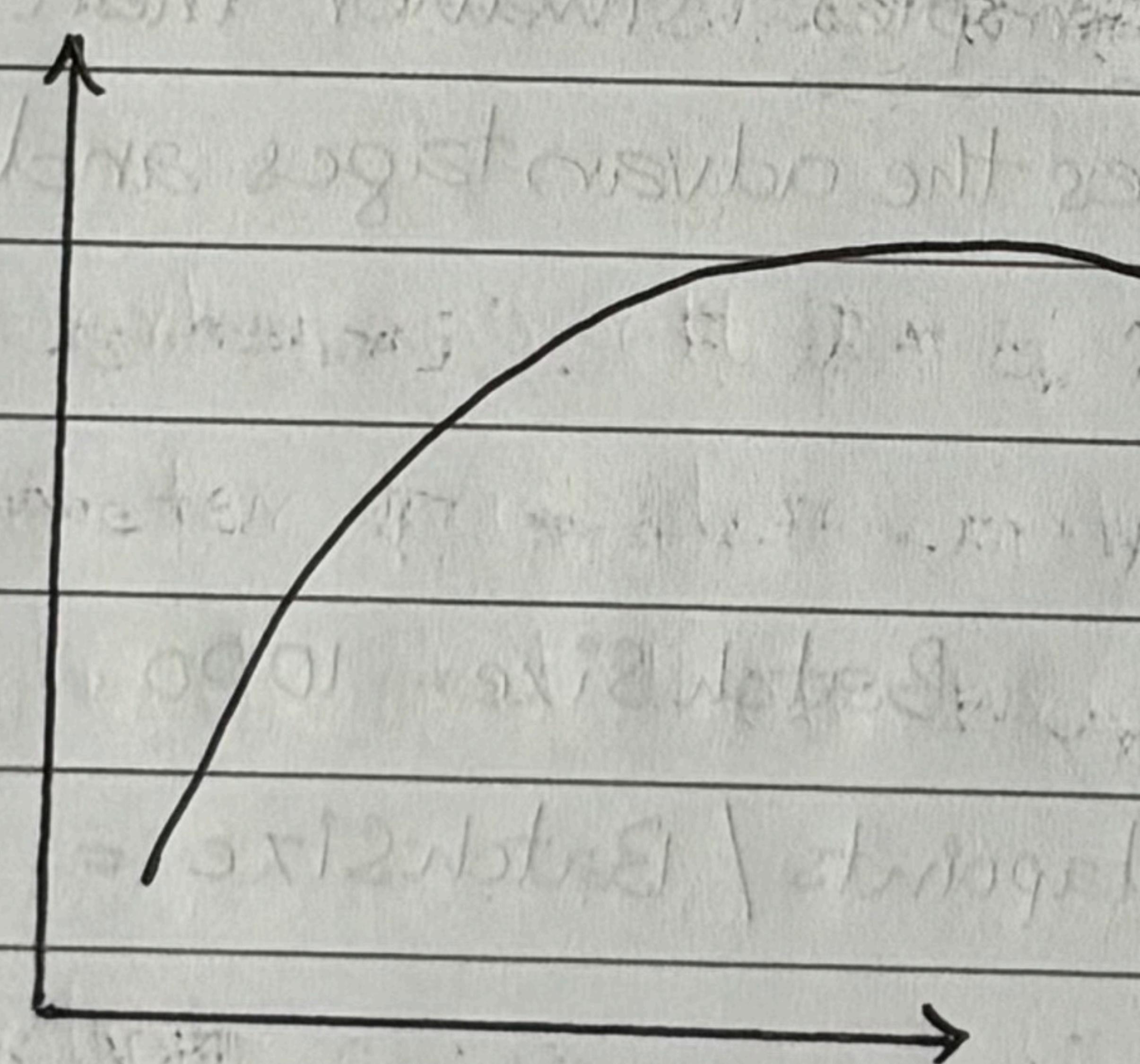
Value $t_1 = a_1, \text{ Value } t_2 = a_2, \dots$



If we introduce a term  $\beta$ , we can smoothen this.

$$\therefore V_{t_2} = \beta * V_{t_1} + (1-\beta) * a_2 \quad \beta = [0-1]$$

If we want to smoothen w.r.t. the previous value give  $\beta \uparrow$  value.



As we can see the curve got smoothened.

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

5. AdaGrad Optimizer: AdaGrad adapts the learning rate for each parameter based on the historical gradient information. The learning rate decreases over time, making AdaGrad effective for sparse features.

→ As the convergence happens the learning rate should change.

$$W_t = W_{t-1} - \eta' \frac{\partial \text{loss}}{\partial W_{t-1}}$$

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

$\eta'$  = dynamic learning rate,  $\epsilon$  = small value

$$\alpha_t = \sum_{t=1}^t \left( \frac{\partial \text{loss}}{\partial W_t} \right)^2$$

As  $\alpha_t \uparrow \rightarrow \eta' \downarrow$

→ Disadvantage:  $\eta'$ , there is a possibility to become a very small value  $\approx 0$

6. AdaDelta and RMS PROPS: Improves upon AdaGrad by introducing a decay factor to prevent the learning rate from decreasing too rapidly.

Here  $\eta' = \frac{\eta}{\sqrt{S_{\text{d}W_t} + \epsilon}}$ ,  $S_{\text{d}W_t}$  → Exponential Weight Average

$$S_{\text{d}W_t} = \beta * S_{\text{d}W_{t-1}} + (1-\beta) \left( \frac{\partial \text{loss}}{\partial W_{t-1}} \right)^2$$

It handles both dynamic learning rate and Smoothening EWA.

7. Adam Optimizer : Adaptive Moment Estimation optimizer combines the advantages of Momentum and RMSProp techniques to adjust learning rates during training. It works well for large datasets and complex models because it uses memory efficiently and adapts the learning rate for each parameter automatically.

$$w_t = w_{t-1} - \eta^t v_{dw}, \quad b_t = b_{t-1} - \eta^t v_{db}$$

→ weight updation

→ bias updation

$$\eta^t = \eta$$

$$\sqrt{s_{dw,t} + \epsilon}$$

dynamic learning

→ smoothening

$$v_{dw,t} = \beta * v_{dw,t-1} + (1-\beta) \delta h_{oss}$$

$\delta w_{t-1} \Rightarrow$  Momentum

$$v_{db,t} = \beta * v_{db,t-1} + (1-\beta) \delta h_{oss} \quad (\text{smoothening weight updates})$$

→ Why Adam works so well ?

1. Each parameter has its own dynamic learning rate based on past gradients and their magnitudes.

2. Bias correction

3. Adam typically requires fewer hyperparameter tuning adjustments compared to other optimization algos making it more convenient choice for most problems.