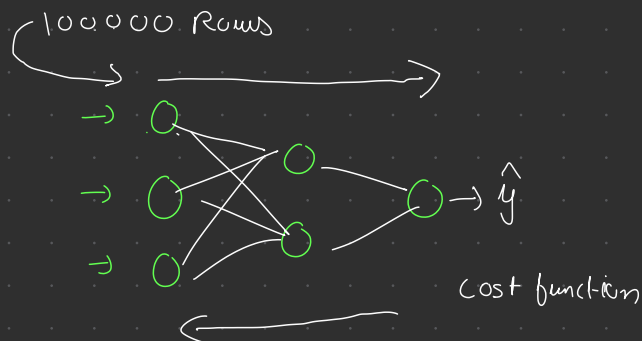


# Optimizers

## 1. Batch Gradient Descent (a.k.a. "Vanilla" Gradient Descent)

- You send all 1 lakh rows at once.
- Compute cKerasfunction on the entire dataset.
- Do 1 weight update per epoch.
- If you run 100 epochs → 100 weight updates in total.

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



### Advantages

- Smooth and stable updates (less noisy).
- Converges steadily because gradient is calculated from all data.

### Disadvantages

- Slow when dataset is huge (computing gradients for all rows before updating).
  - Requires a lot of RAM/VRAM to load the full dataset at once.
- 100 epochs → updating the weights 100 times

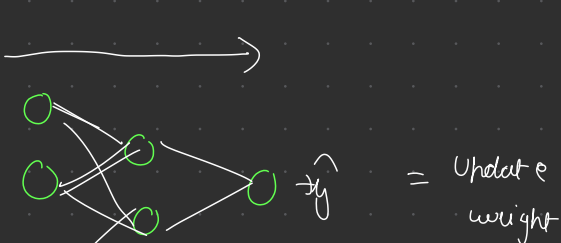
## 2. Stochastic Gradient Descent (SGD)

- You send 1 row at a time.
- Compute cost for that single row.
- Update weights immediately.
- With 1 lakh rows, in 1 epoch → 1 lakh updates.
- In 100 epochs → 1 crore updates.

100000 Rows

100 epoch

100 x 100000  
= 1 crore UP



### Advantages

- Much faster to start learning (weights update after every row).
- Works well for very large datasets (you don't need all rows in memory at once).

### Disadvantages

- Updates are very noisy → loss curve jumps around instead of smoothly decreasing.
- Can be unstable (harder to converge).
- Slower to reach the exact minimum compared to mini-batch.

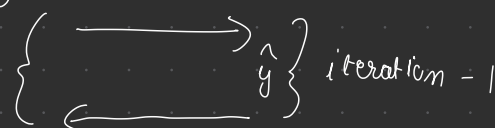
2 Row  
3 Row  
4 Row  
⋮  
100000 Rows

## 3. The Practical Solution: Mini-Batch Gradient Descent

- A compromise between batch and SGD.
- Split data into small batches (e.g., 32, 64, 128 rows).
- Each batch → forward pass → backward pass → weight update.
- So: in 1 lakh rows, batch size 100 → 1000 updates per epoch.

100 epoch

100000 ⇒ 100 Batch size

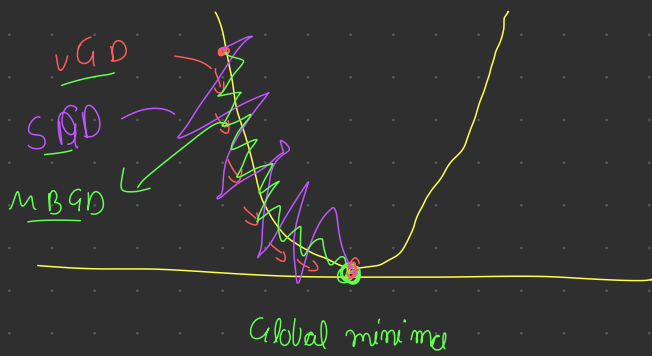


This is what almost all deep learning frameworks use today (including Keras/TensorFlow).

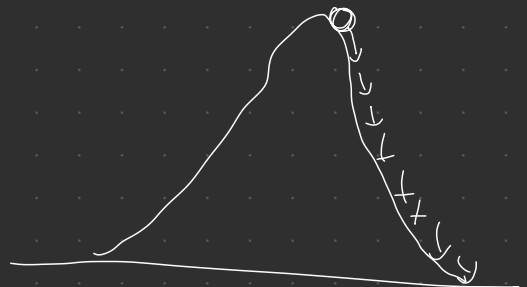
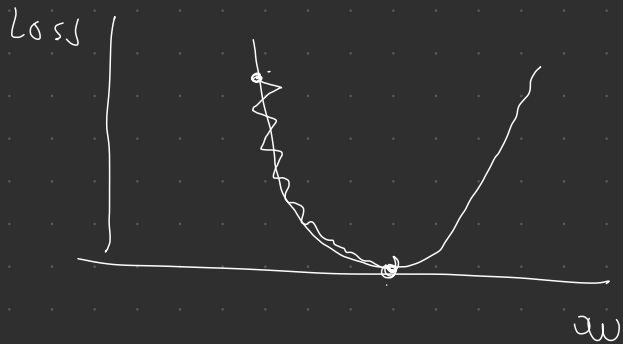
$$100 \times 1000 = 100000$$

100000  
100

1000



Momentum with SGD



Momentum =  $\boxed{\text{mass}}$  x velocity

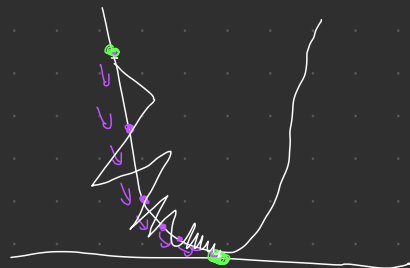
$$w_{\text{new}} = w_{\text{old}} - \eta \left[ \frac{\partial L}{\partial w_{\text{old}}} \right] \Rightarrow w_{T+1} = w_T - \eta \Delta w$$

$$w_{T+1} = w_T - \boxed{v_t} \quad \hookrightarrow \quad v_t = \beta \times v_{t-1} + \eta \Delta w_t$$

$\boxed{0.09}$

Adagrad (Adaptive Gradient descent)

$$w_T = w_{T-1} - \underbrace{\boxed{\eta}}_{\substack{0.01 \\ \text{constant}}} \frac{\partial L}{\partial w_{T-1}} \longrightarrow$$

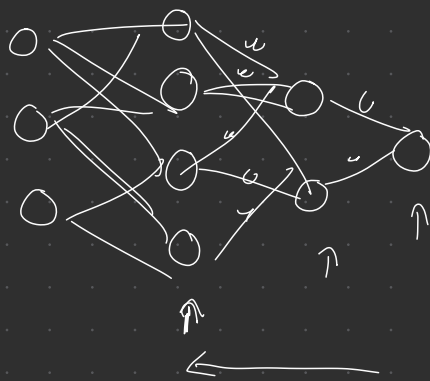


$$w_T = w_{T-1} - \boxed{\eta} \frac{\partial L}{\partial w_{T-1}}$$

$$\Downarrow \eta' = \eta$$

$$\rightarrow \boxed{\alpha_T + \epsilon} \quad \text{Epsilon}$$

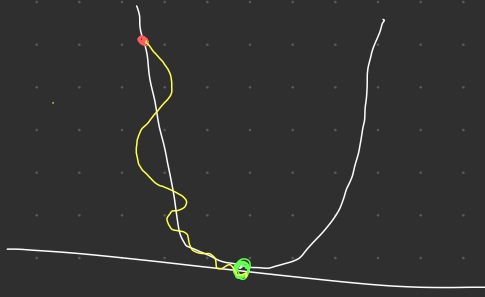
$$\alpha_T = \frac{1}{\sum_{i=1}^T} \left( \frac{\partial L}{\partial w_T} \right)^2 \Rightarrow$$



$$0.01 \rightarrow 0.002 \rightarrow \underline{0.001}$$

Adam optimizer

↳ Adagrad + Momentum



Decision Tree  $\rightarrow$  white box model

Random Forest  $\rightarrow$  Black box model

ANN  $\rightarrow$  Black box model

LR - white box model