

Optimizers

What are Optimizers ?

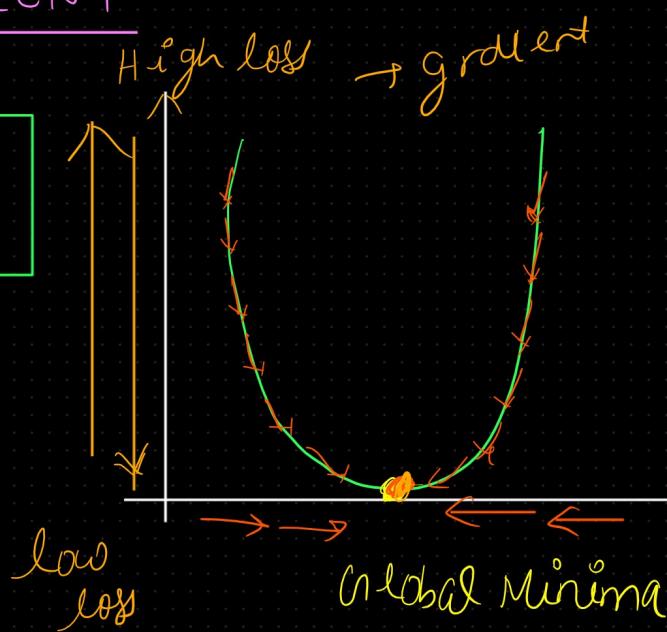
Optimizers are algorithms used in training deep learning models to minimize the loss function and improve the model's performance.

Some common optimizer functions :-

- (1) Gradient Descent (GD)
- (2) Stochastic Gradient Descent (SGD)
- (3) Min - Batch Gradient Descent
- (4) Momentum
- (5) Adagrad
- (6) RMSprop
- (7) Adam (Adaptive Moment Estimation)

(1) GRADIENT DESCENT

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

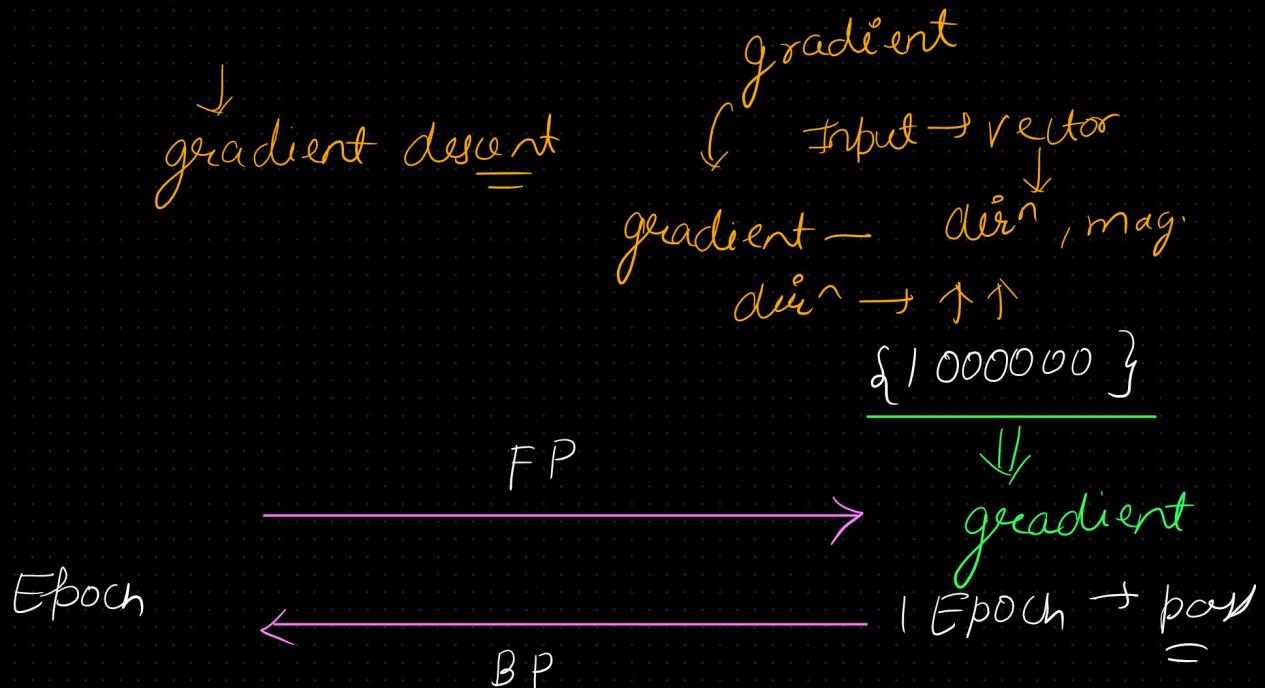


GD → GD is a first-order iterative optimization algorithm for finding the minimum of a function.

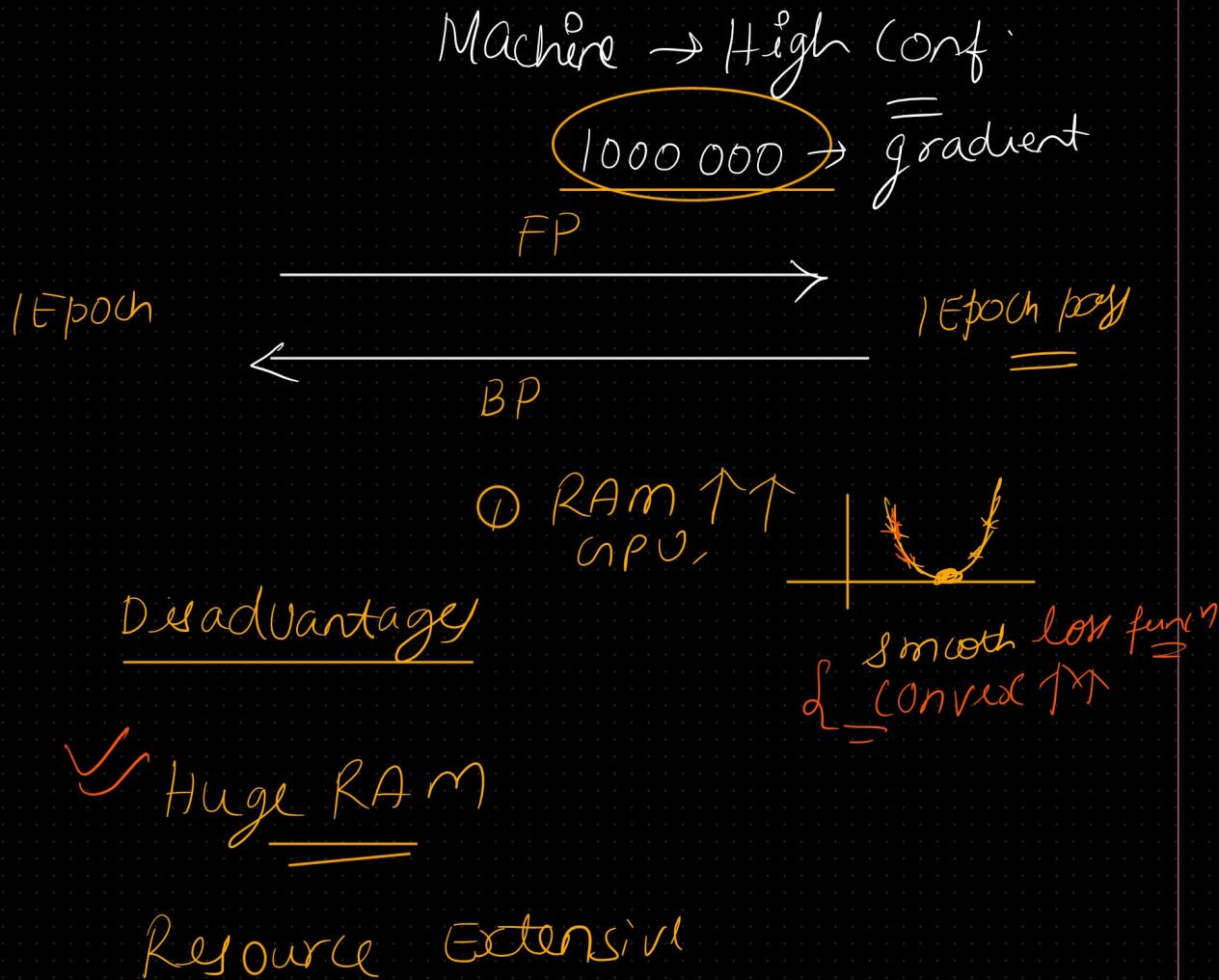
Intuition

It works by iteratively moving in the direction of the negative gradient of the loss function with respect to the model parameters.

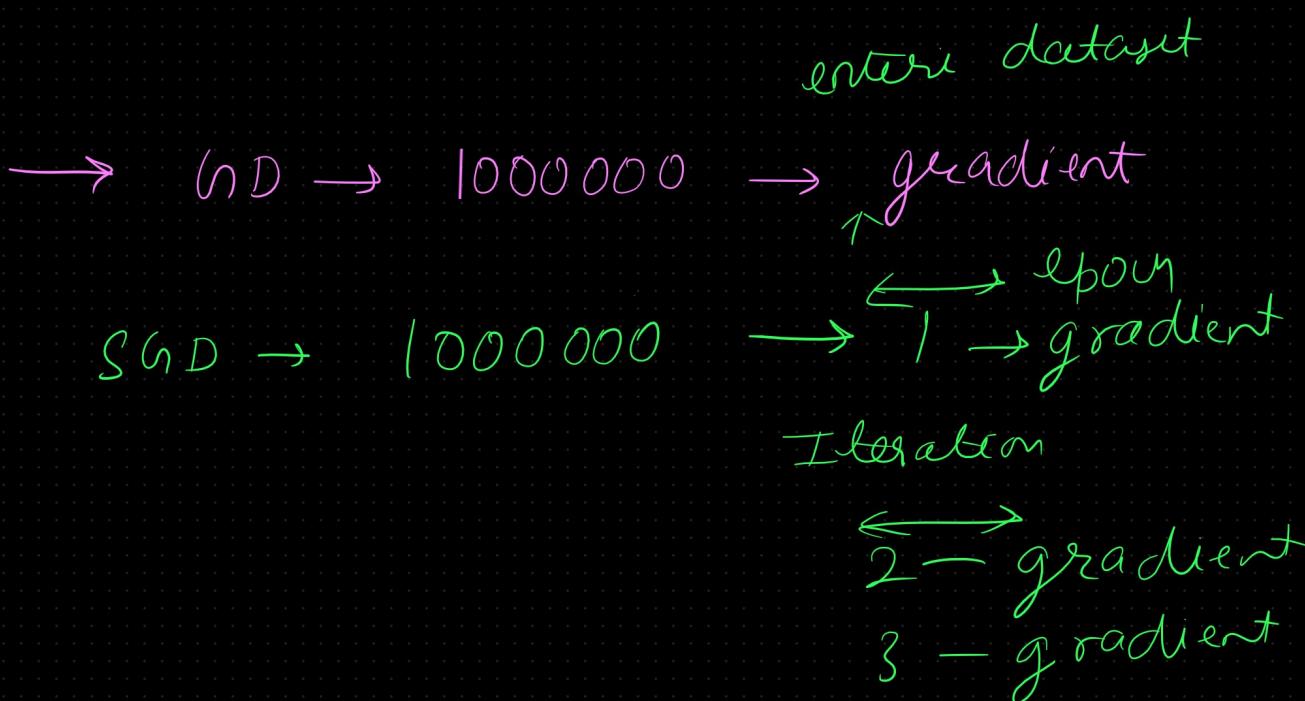
(weight, bias)



Weight update
Minimum loss function



② Stochastic Gradient Descent (SGD)



1000 000

1 Epoch

1 record $\xrightarrow{\text{weight update}} \hat{y}^{\text{1st}}$

2 record $\xleftarrow{\text{weight update}} \hat{y}^{\text{2nd}}$

⋮
⋮

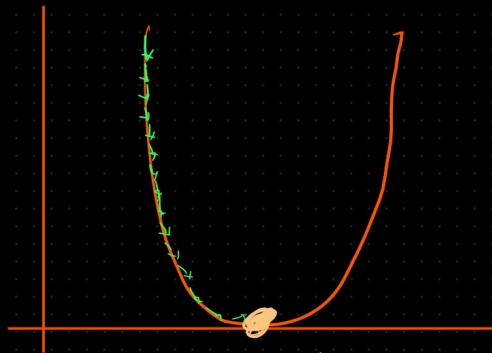
$\xleftarrow{\text{weight update}} \hat{y}^{\text{million}}$

SGD is an extension of GD where instead of computing the gradient of the entire dataset, it computes the gradient of the loss for each training example and updates the parameters accordingly.

Advantage

1) RAM \downarrow

2) More Computationally Efficient than GD



Disadvantage

1) Convergence will be very slow

2) Time Complexity $\uparrow\uparrow$

3) Mini Batch SGD

1000000 → DATASET

↓

1000000

1000

Batch = 1000
size

= 1000
Iteration

1000
records

→ } Iteration
1

1000
records

→ } Iteration
2

1000
records

→ } Iteration
m

Mini-batch GD is a compromise between GD and SGD. It divides the dataset into small batches and computes the gradient of the loss for each batch.

Advantage

- ① Faster convergence than SGD

Disadvantage

- ① Resource Intensive
- ② Require tuning the Batch size which affect the convergence and memory usage

(4) Momentum \rightarrow Smooth

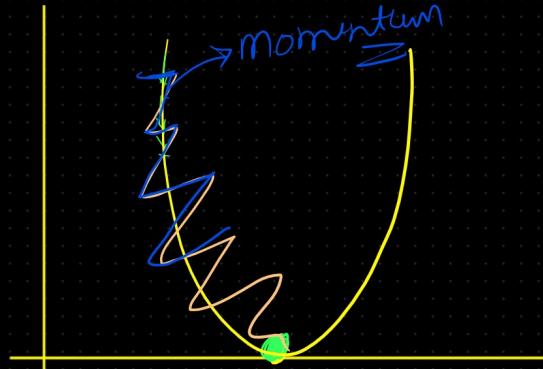
time series data
=

Weight updation

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

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

Momentum is a technique that accelerates SGD by adding a fraction of the update vector of the past time step to the current update vector.



It helps to smooth out variations in the gradient descent path, especially in the presence of noisy gradients or high curvature.

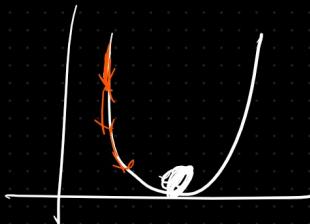
(5) Adagrad

↙

α - learning Rate \rightarrow Adaptive

Adagrad adapts the learning rate of each parameter based on the historical gradients for that parameter.

Intuition



$$\frac{\partial l}{\partial w} = \text{gradient} \approx \alpha \rightarrow \text{scale down}$$

$$\text{gradient} \downarrow \approx \alpha - \text{scale up}$$

Disadvantage

Accumulates the squared gradients in the denominator, leading to diminishing learning rates over time.

(6) RMSprop

RMSprop is an adaptive learning rate optimization algorithm that divides the learning rate by an exponentially decaying average of squared gradients.

$$n' = \frac{n}{\overline{\sigma}}$$

$$n' = \frac{n}{\sqrt{\sigma_{t+1} + \epsilon}} \quad \text{Exponential}$$

→ It scales the learning rate differently for each parameter based on the magnitudes of recent gradients.

↓ Best (7) Adam (Adaptive Moment Estimation)

→ Adam combines the ideas of momentum and RMSprop. It computes adaptive learning rates for each parameter and keeps track of both the first and second moments of the gradients.

→ It is less sensitive to the choice of hyperparameters and can handle sparse gradients more effectively.

Advantage

- 1) Fast Convergence
- 2) Good Generalisation

Disadvantages

Introduces additional hyperparameters (momentum coefficients and decay rates) that need to be tuned.

May exhibit erratic behavior on certain types of problems or architectures if hyperparameters are not properly adjusted.