

Optimizers in Neural Network

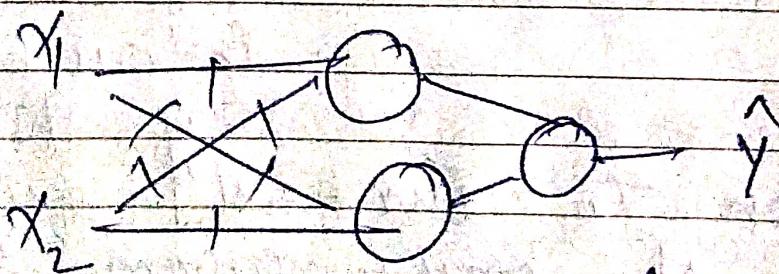
While performing any neural network model training we're always concerned of the following things

- 1) Performance of ANN
- 2) How to speed up the training?

3 techniques to speed up:

- 1) weight initialization
- 2) Batch Normalization
- 3) In some cases Activation fn.
- 4) Optimizers.

What actually optimizers do?



4 2 2 1 (bias)

weights ↓ ↓ q parameters
bias weights

We need to find such a value from the q parameters that whatever the input model gives should be close to

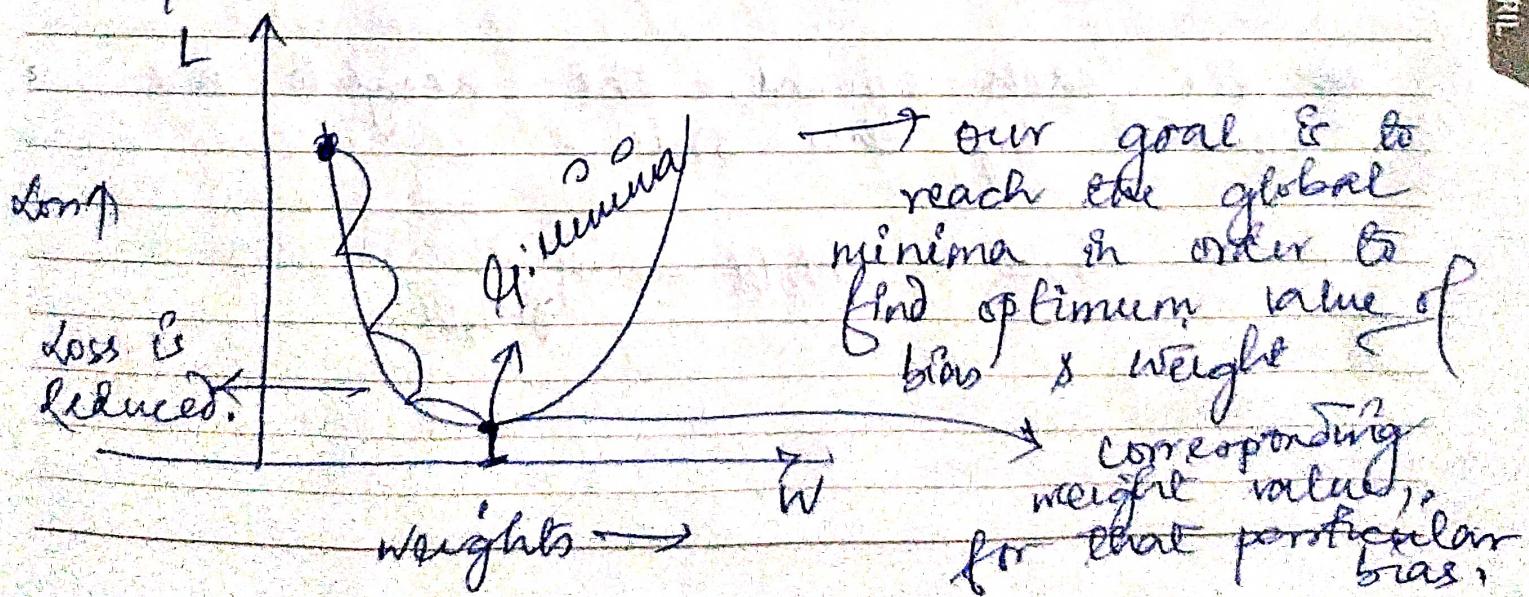
the original value.

→ Thus, we can conclude that the neural network is actually facing the optimisation problem wherein,

- (1) we want to reduce the loss.
- (2) And to achieve this, the set values of weights & biases will be the optimum parameter of for our ANN.

for this, we assign random values to weights & bias, & slowly & gradually while training, the optimum value of the weights & biases are reached.

If we draw the graph in 2D, as for a parameter, → 1D graph is plotted.



- So, to perform this task, we use Gradient Descent.
- In DL, we use Gradient Descent, Gradient Descent,

$$w_1 = w_0 - \eta \frac{\partial L}{\partial w}$$

↓ new weight ↓ old weight learning rate

→ 3rd Type of Gradient Descent:-

- (I) Batch GD
- (II) Stochastic GD.
- (III) Mini-Batch GD.

→ For each epoch, the weights are updated by this rule.

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

→ The main difference b/w the mentioned 3 AD, is that,

✓ How many datapoints, you want to see, before updating the weight!

Difference factor

(i) for Batch AD, the whole datapoints are traversed.

weight is updated according to the epoch. For each epoch, weight updation is done, (for 10 epochs, 10 times the weight will be added).

(ii) for stochastic AD; the entire datapoints are traversed but for each datapoint, there will be weight updation.

For eg, If there are 500 datapoints,
then 500×10 (epochs)

≥ 5000 times
weight updation takes place.

(iii) In mini-batch, we need to decide a batch for the available datasets & divide into batch for model training.

For example, for 500 datapoints, a batch of 100 datapoints may be used for training the model architecture.

5 batch \rightarrow 10x5 times

↑
epochs

→ Challenges with previous optimizers:

→ Decide learning rate

$$w_n = w_0 - \eta \frac{\partial L}{\partial w_0}$$

→ If we take an ideal case → if we take a very less, convergence will be very slow.

→ Also, if we take a bigger value, convergence will be very fast, may shoot the global minima.

bigger η

Solⁿ: 2) Learning Rate scheduler:

- pre-defined training model problem of this, when to increase or decrease the learning rate as it is (after how many epochs) defined before the model training.
- It means it doesn't depend on any dataset. As if it is dependent on dataset, to bring out proper result would be challenging.

3) We need to minimize multiple weights.

Ex: If there are two weights w_1 & w_2 , we cannot assign separate learning rate for both the params, gradient descent doesn't provide that flexibility.

4) Local minima:

- The GD graph is bit complex, our target is to reach the global minima, but we fall into the problem of local minima, & get stuck & cannot get out of that.

SGD has that ability to get out

as its movement is like zig-zag, but for others it is quite not possible, & work with auto-optimization methods, soh.

5) Saddle point:-

A point where slope of one (say) dimension moves in one direction, some other slope moves in other directions, (say) this creates a place where, the point has the slope doesn't change, i.e., the slope remains constant. ($\text{diff} = 0$), so, while we're updating, the word = when, so, we get stuck.

for this reason, we need to use some different optimizers:-

- 1) Momentum
- 2) Adagrad
- 3) NAG
- 4) RMS Prop
- 5) Adam

→ Exponentially weighted Moving Average.