

Optimization in Deep Learning

- In Deep Learning, with the help of loss function, the performance of the model is evaluated.
- This loss is used to train the network so that it performs better.
- Essentially, we try to minimize the loss function.
- Lower loss means the model performs better.
- The process of minimizing any mathematical function is called optimization.
- Optimizers are algorithms used to change the features of neural network such as weight and learning rate so that the loss is reduced.
- The goal of an optimizer is to minimize the objective function.

Need for Optimization

- ~~Presence~~ Presence of local minima reduces the model performance.
- To minimize the loss value (Training error).
- To select appropriate weight values and other associated model parameters.

Types of optimization

1. Gradient descent

- it starts with some coefficients seen
- it moves towards lower weight and updates the values of coefficient and repeat until the local min is reached

Disadvantages

- Expensive to calculate a gradient if the size of the data is huge
- Not suitable for non-convex function

2. Stochastic Gradient Descent :

- Instead of taking the whole data set for each iteration randomly select batches of the data
- Select the initial parameter w and learning rate
- Randomly shuffle the data and each iteration to reach the approximate minimum.

- fast is than GD but the computation
- Since only few batches are

Disadvantage 1

3. Stochastic Gradient descent with Momentum:

- Since SGD is a noisy path we are going for SGD with momentum
- momentum helps in fast convergence of the loss function
- SGD oscillates b/w either direction of the gradient and update the weight.
- By adding the fraction of the previous update to the current update will make the process a bit faster.

4. Mini batch gradient descent:

- Only a subset of the dataset is used for calculating the loss function
- It takes only fewer iterations so faster than SGD
- It is smoother than SGD
- It has good balance b/w speed and accuracy

5. Adagrad Adaptive gradient descent

- It uses different learning rates for each iteration
- The change in the learning rate depends upon the difference in the parameter unit training.

and