

Epoch, Iteration

13 September 2023

07:07

Loss vs steps vs epochs

→ 1000 → 1000

Epoch

Whenever entire dataset is being passed to Neural Net is called as an epoch.

10M

10M → Data → NN → \hat{y} → Loss ← WB. 1 Epoch

10M → Data → NN → \hat{y} → Loss ← WB. 2nd Epoch.

Iteration

Whenever the subset of data set is being passed to Neural Net is called as iteration.

1000R

100 → NN → \hat{y} → Loss ← WB. Iteration 1

100 → NN → \hat{y} → Loss ← WB. Iteration 2

Iteration 10th

1000
100 100 100
→ complete 10 Iteration } 1 Epoch

1 Epoch could consist of multiple iteration

Model

1000R → 5 Epoch. 5 Iteration.

25 Iteration.

100 Random Data → NN → \hat{y} → Loss ← WB. Loss = 100.6

WB 100 Data → NN → \hat{y} → Loss ← WB

WB 100

Loss ← 10th. 1 Epoch. } Epoch.

9D was only the algorithm that works on Epoch concept

Stochastic Gradient Descent

Random.

In case of gradient descent we use to take entire dataset in forward propagation

	x_1	x_2	x_3	x_4	Yact
✓ R_1	^	^	^	^	^
R_2	^	^	^	^	^
✓ R_3	^	^	^	^	^
R_4	^	^	^	^	^

Random

R_0 → NN → \hat{y} → Loss ← WB

R_1 → NN → \hat{y} → Loss ← WB

R_4 → NN → \hat{y} → Loss ← WB

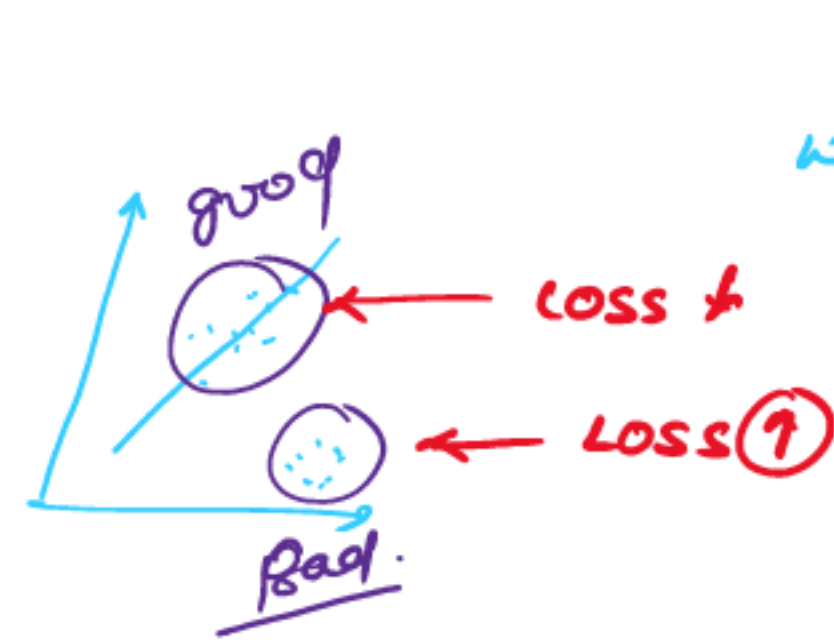
R_2 → NN → \hat{y} → Loss ← WB

4 Iteration 4 time WB

1 Epoch

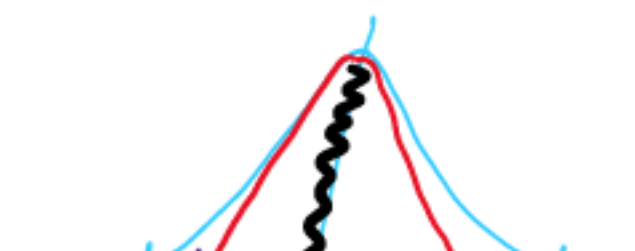
1M SGD 1 Epoch 1WB

SGD 1 Epoch 1MWB

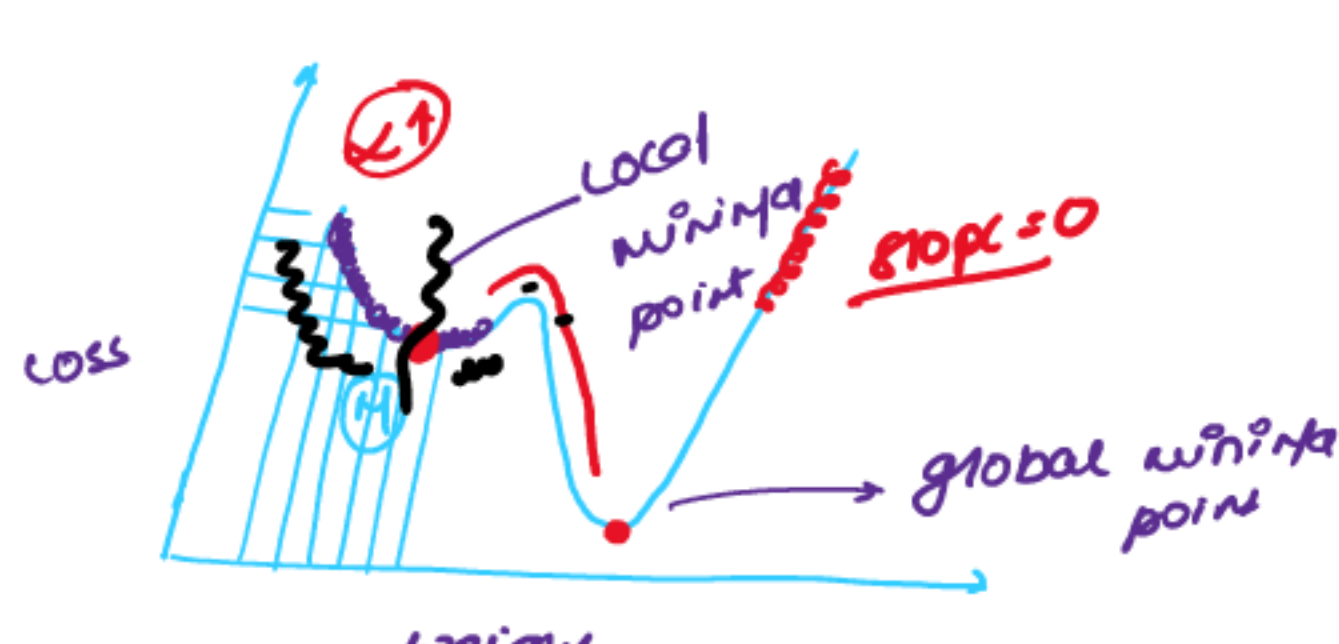


$$w_{new} = w_{old} - \alpha \frac{\partial L}{\partial w_{old}}$$

$$0.1 - 0.01 = 0.09$$



SGD is prone to multiple local minima.



$R_3 = \text{good} = \text{Loss} \downarrow$

$R_1 = \text{good} = \text{Loss} \downarrow$

$R_4 = \text{good} = \text{Loss} \downarrow$

$R_2 = \text{Bad} = \text{Loss} \uparrow$

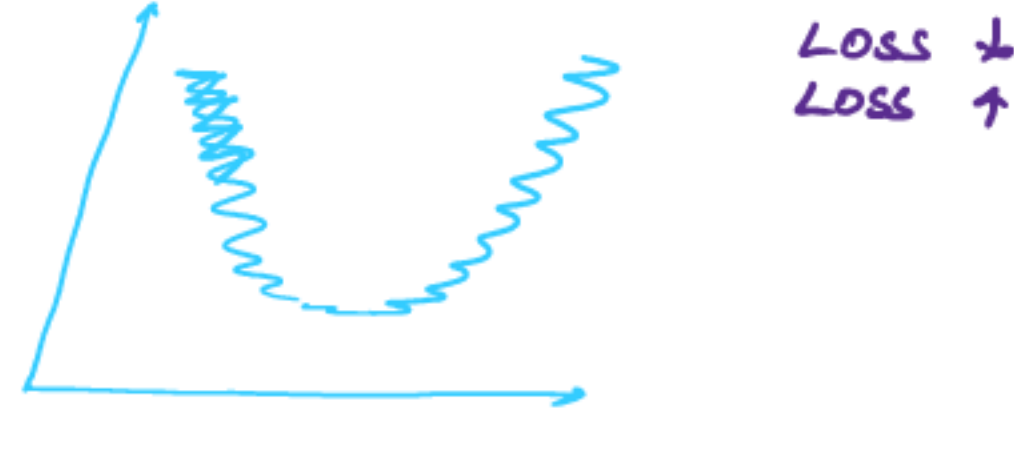
It is selecting random data points from dataset and calculating derivative as we can say updating the parameters

Yact $\hat{y} = \text{sinusoid}$ WB

Yact $\hat{y} = \text{Loss}$ WB

* Batch size = 1
no. of Batches = Data size

path is noisy
Computationally not expensive.



Loss ↑ WB ↑

Loss ↓ WB ↓

$$\frac{y_{act} - \hat{y}}{15 - 14} = \frac{1}{15}$$

$$\frac{y_{act} - \hat{y}}{15 - 5} = \frac{1}{10}$$

Test 8P

as on sample from data is chosen randomly path taken to

Reach global minima is usually faster than gradient descent.



Derivative

Iteration 10