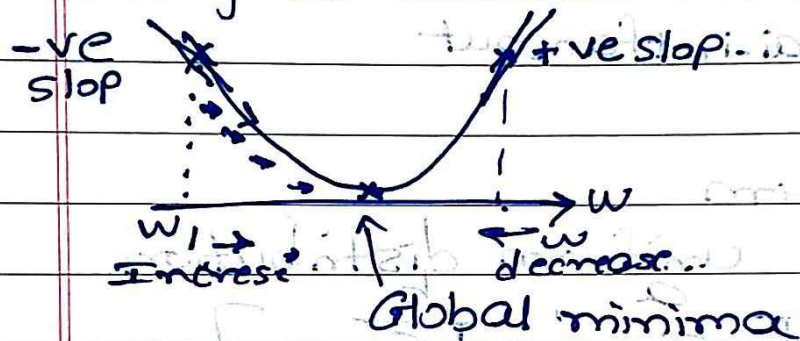


## \* Stochastic Gradient Descent

In backpropagation weight updation formula<sup>is</sup>

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial L}{\partial w_{\text{old}}} \quad (1)$$

In gradient descent



In gr

1) In eqn (1) for weight updation consider all training example for one epoch & update weight by backpropagation. This is a gradient descent algorithm.

2) Where as in stochastic gradient descent consider one training example in one epoch & update the weight.

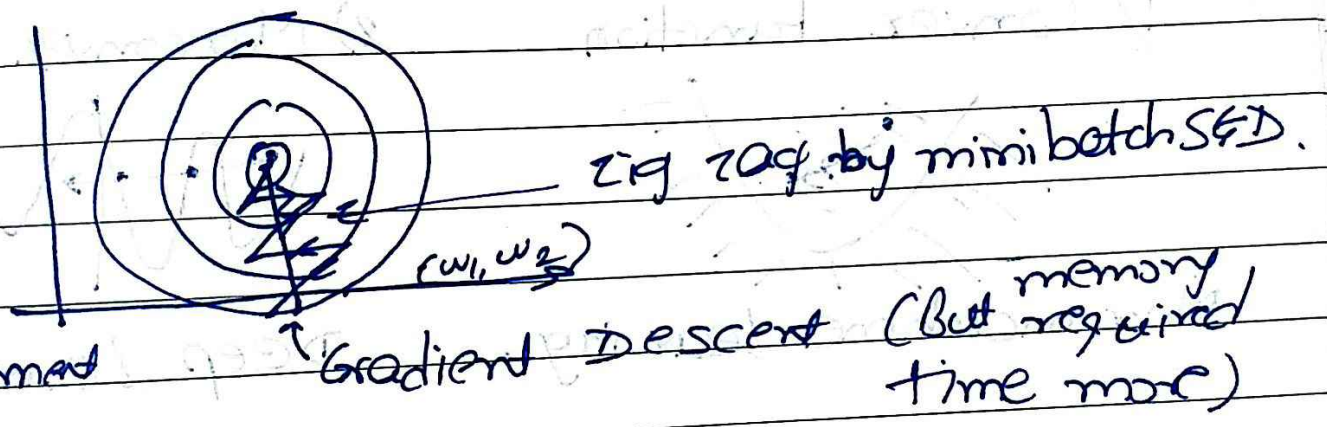
This is SGD or minibatch SGD



In many ANN & CNN minibatch SGD used.

Because when we have million of record so to load this much record it take lot computation time

In order to prevent this mini batch SGD used.



\* So in mini batch SGD we use K datapoint as minibatch zigzag occurs so to avoid that momentum used.

$$L(w) = \frac{1}{K} \sum_{i=1}^K (y - \hat{y})^2$$

K = Batch of training example

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

so to avoid computation time use SGD

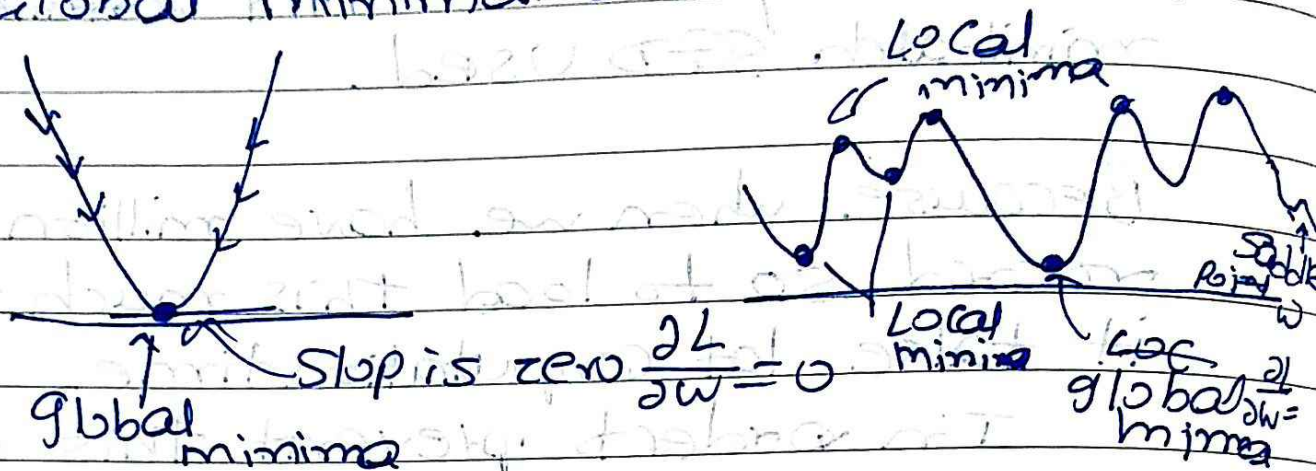
$$\left[ \frac{\partial L}{\partial w_{old}} \right] \approx \left[ \frac{\partial L}{\partial w_{old}} \right]_{GD}$$

↑ population

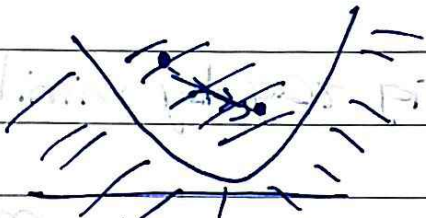
mini batch SGD  
sample



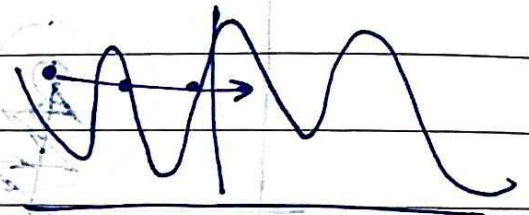
# Global minima Vs Local minima



1) Convex function



2) Nonconvex function



Normal machine learning

Deep Learning

In deep learning as so many layer & neuron are present due to which nonconvex function having local & global minima.

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

$$= w_{\text{old}} - \left[ \eta v_{t-1} + \eta \frac{\partial L}{\partial w_{\text{old}}} \right]$$

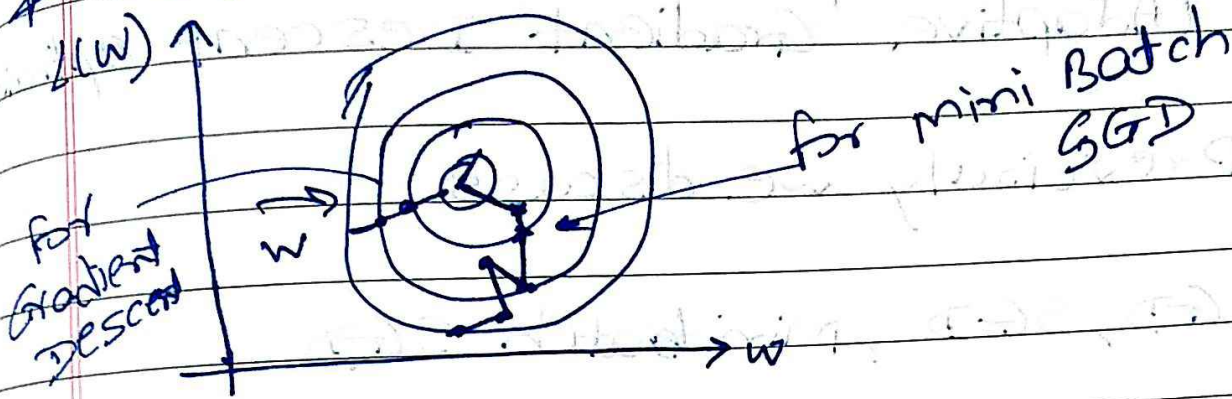
$$v_{t-1} = \eta \times \left[ \frac{\partial L}{\partial w_{\text{old}}} \right]_t + \eta^2 \left[ \frac{\partial L}{\partial w_{\text{old}}} \right]_{t-2} + \dots$$

Data point at different time interval.

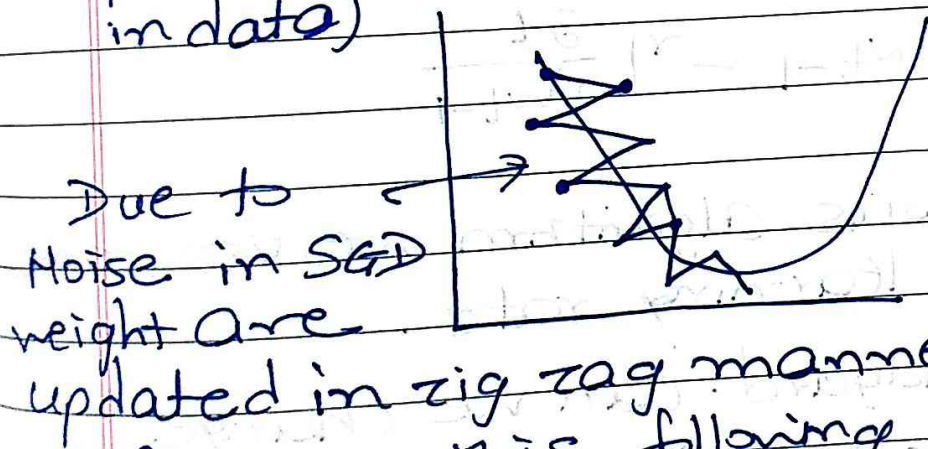
$$\eta = 0 \text{ to } 1$$



## \* SGD With Momentum



1) In minibatch SGD we take  $k$  data point (Some times noise gets added in data)



Due to Noise in SGD weight are

updated in zig zag manner

Because this following problem occur.

- 1) for convergence it take more time to reach to global minima because batch
- 2) for the same we need to reduce noise

So to reduce noise use ↴

\* Exponentially moving Average  
ie. Momentum



# \* Adgard Optimizer (Adaptive Gradient Descent Optimizer) \*

Previously we discuss:

1) GD, SGD, minibatch SGD

2) Consider weight updation  

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

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

3) In previous algorithm we use same weights & learning rate.

But in Adgard. can we change different learning rate for different iteration number.

4) we consider different learning rate because of

1) Dense feature (f)

2) Sparse feature  $\leftarrow$  BoW (many are zero, some one)

1) In dense most of the value are non zero

2) In sparse most of the value are zero



are zero.

So in Adagrad due to dense & sparse we use different learning rate using on different iteration number.

$$w_t = w_{t-1} - \eta_t \times \frac{\partial L}{\partial w_{t-1}} \quad \text{t: No of iteration}$$

$$\eta_t = \frac{\eta}{\sqrt{t + \epsilon}} \quad \eta = \text{initial learning rate}$$

$$\alpha_t = \sum_{i=1}^t \left( \frac{\partial L}{\partial w_{t-1}} \right)^2 \quad \eta_t = \frac{\eta}{\sqrt{\alpha_t + \epsilon}}$$

↑ small +ve number.

As  $\alpha_t$  increases  $\eta_t$  decreases as iteration goes on increasing, due to which weight decreasing slowly.

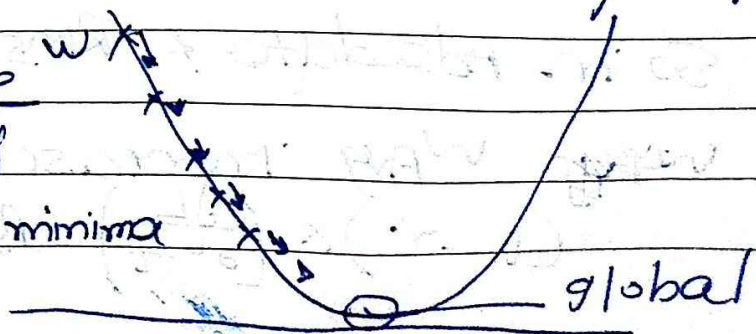
So adagrad use different different learning rate by which we can update weight & optimize global minima.

Limitation of Adagrad ← As sometime iteration increases  $\alpha$  increases

Some time  $\alpha_t$  become very high number

It will handle dense & sparse feature &

Converges for global minima





## \* Adadelta & RMSprop optimizer

In Adagrad

$$W_t = W_{t-1}$$

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

$$\eta'_t = \frac{\eta}{\sqrt{\sum_{i=1}^t \epsilon}}$$

$$\alpha_t = \sum_{i=1}^t \left( \frac{\partial L}{\partial w_{t-1}} \right)^2$$

So here  $\alpha_t$  is high due to which  $\eta'_t$  becomes very small & weight updation not happen.

So in Adadelta

$$\eta'_t = \frac{\eta}{\sqrt{W_{\text{Avg}} + \epsilon}}$$

weighted avg.

Exponential avg.

$$W_{\text{Avg}} = r W_{\text{Avg}, t-1} + (1-r) \left( \frac{\partial L}{\partial w_t} \right)^2$$

$$\alpha_t = \sum_{i=1}^t \left( \frac{\partial L}{\partial w_i} \right)^2 \quad r = 0.95 \text{ selected}$$

So in Adadelta & RMS. we rescaled

~~W Avg~~  $W_{\text{Avg}}$  increase by adding  $(1-r) \times \left( \frac{\partial L}{\partial w_t} \right)^2$  By selecting  $r = 0.95$   
 $\therefore 1-r = 0.05$

classmate

Date \_\_\_\_\_

Page \_\_\_\_\_

so by restricting this learning rate  
also changes slowly.  
& Converges take less time.