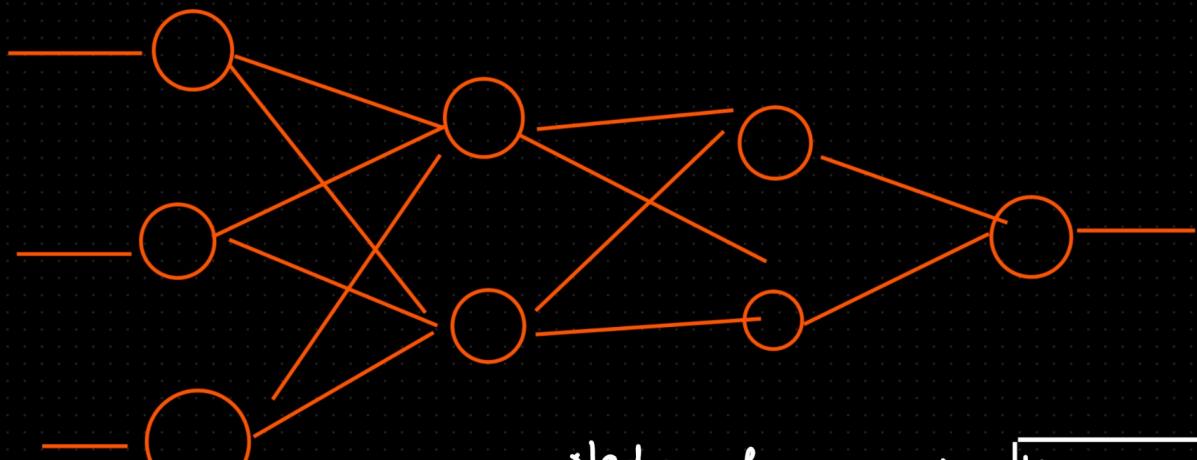
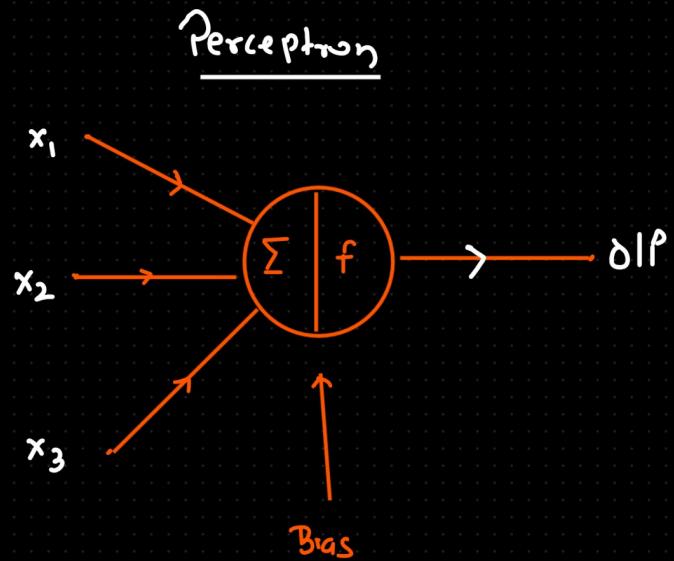


## Optimizers

Gradient-Descent-



Node, Layers  $\Rightarrow$  Hyperparameters

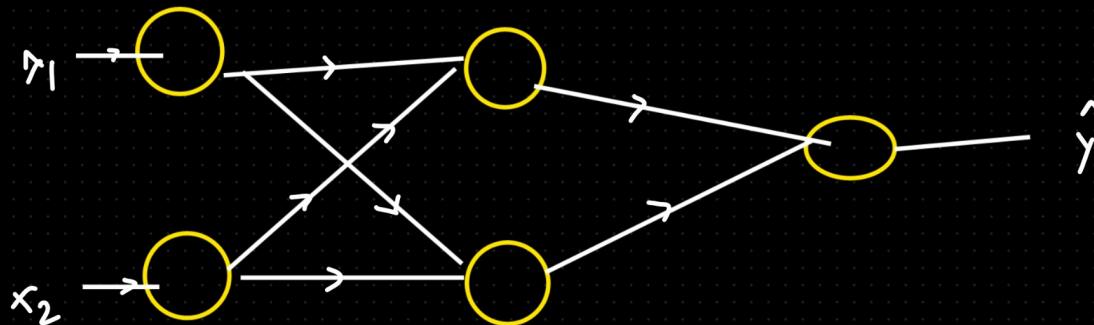
= Complex Neural Network

Performance of Neural Network might Decrease

training time will be increased

? >>

- 1 Optimized weight initialized ✓
  - 2 Batch Normalization ✓
  - 3 Right selection of Algo Alg
  - 4 Optimizer → Right Optimizer
- Performance of Deep NN
- = Decrease the training time



$$2 \times 2 + 2 + 2 \times 1 + 1 = 9$$

6 weight and 3 bias

NN

Loss !!!

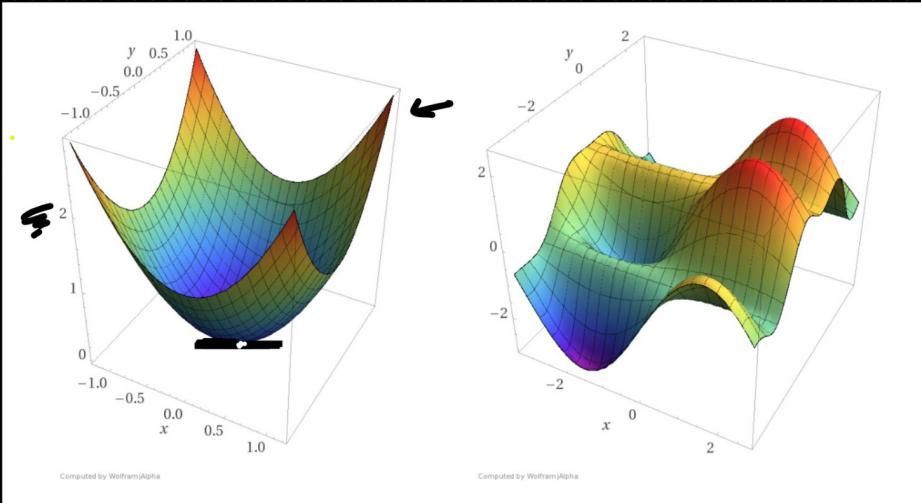
$$\mathcal{L}(\text{weight, bias})$$

$$(y - \hat{y})^2$$

$$y = f(x)$$

$$y = x^2$$

GD



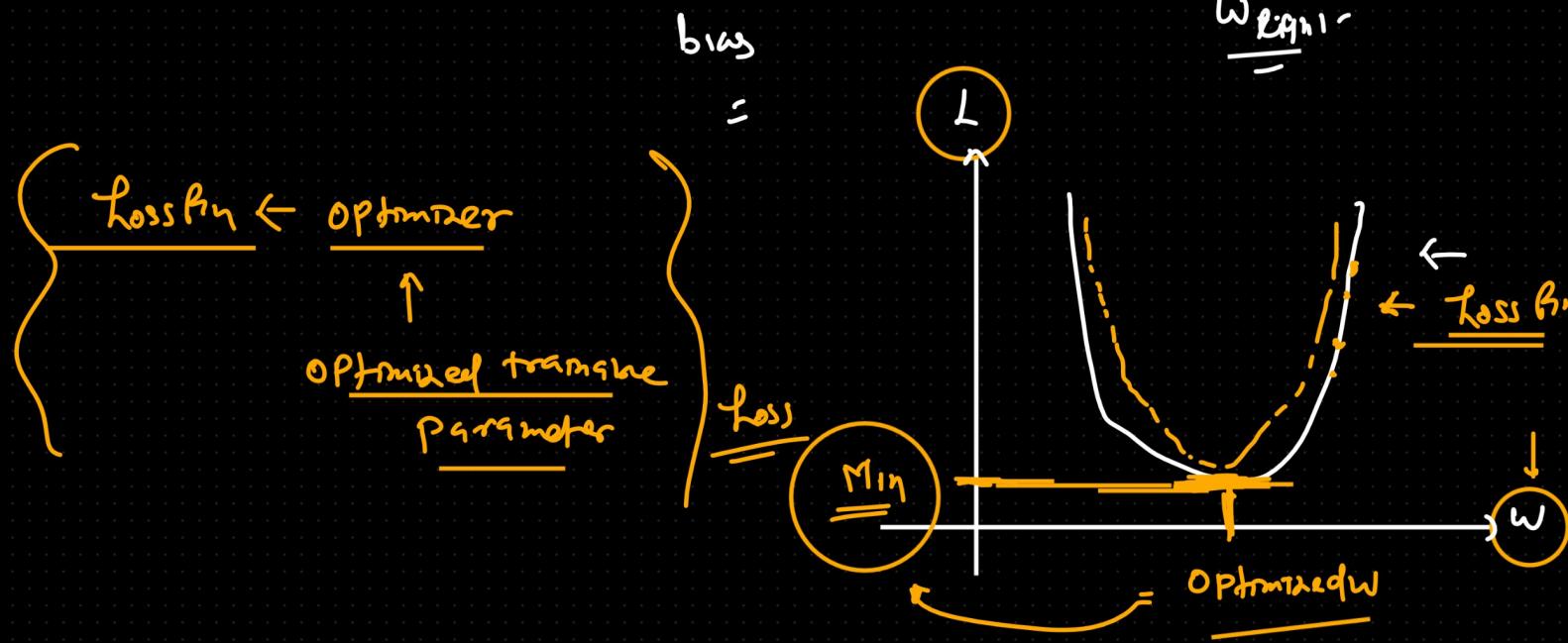


L.D

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

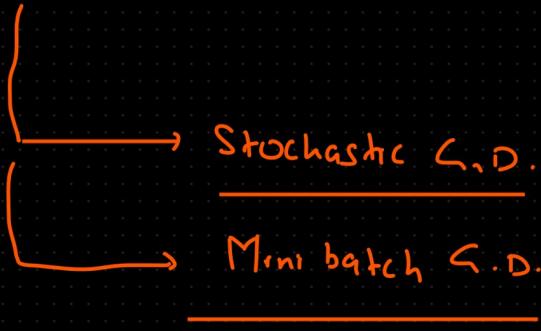
G.D

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



Optimizer  $\Rightarrow$  to optimize the trainable parameter  $\Rightarrow$  minimum loss

= ① Gradient Decent  $\Rightarrow$  Vanilla GD or batch GD.

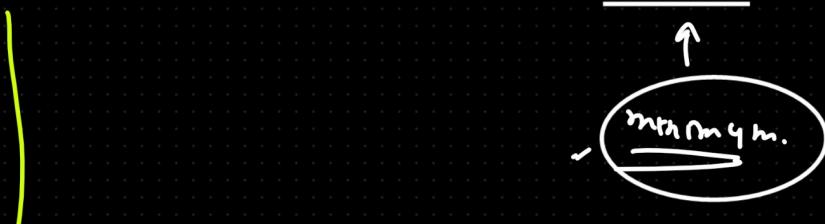


Gradient Decent  
(Vanilla G.D.)

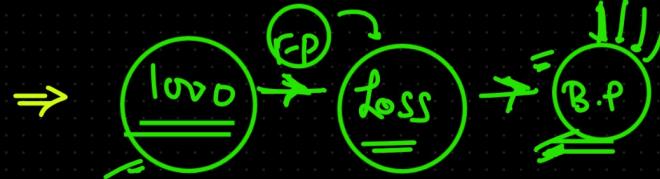
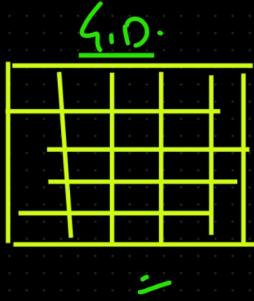
Stochastic GD

Minibatch G.D.

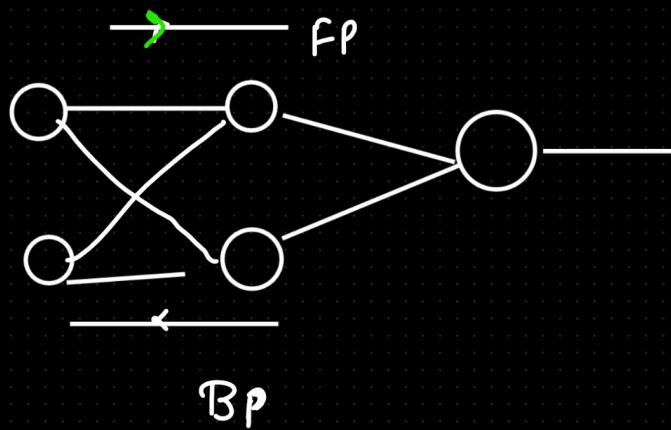
Data  $\Rightarrow$  NN  $\Rightarrow$  OIP  $\Rightarrow$  Loss( $y - \hat{y}$ )  $\Rightarrow$  BP  $\neq$  FP



Number of Rows  $\Rightarrow$



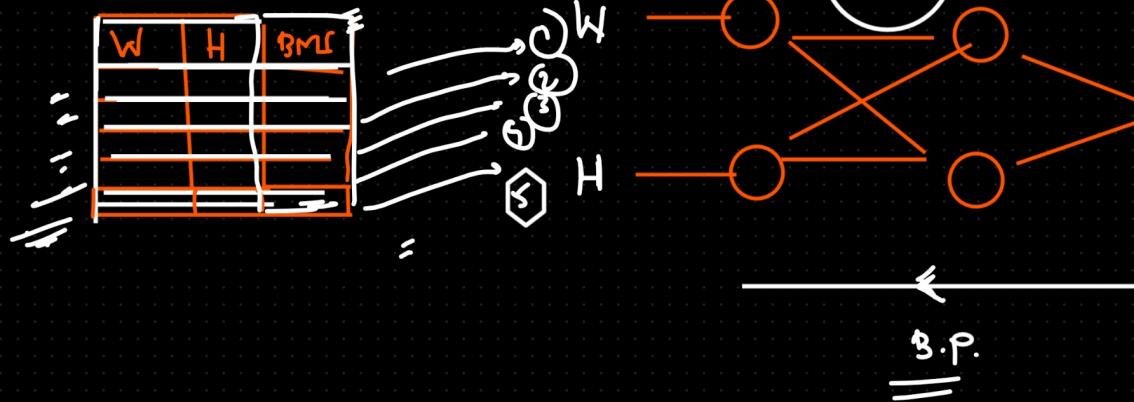
Epoch  $\Rightarrow$  1 forward Propogation + 1 Backward Propogation



10 times

1 epoch

Vanilla G.D

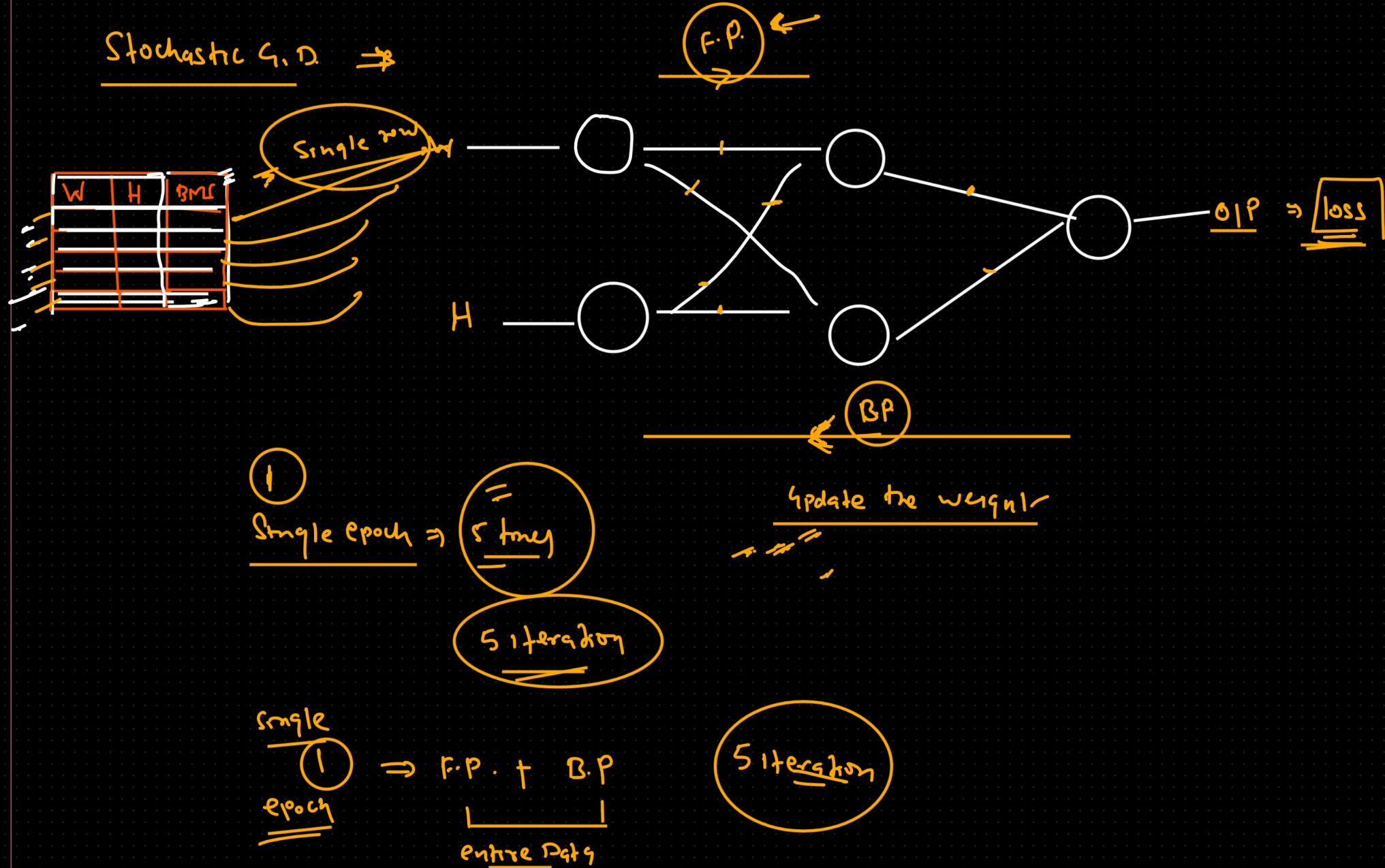


$$\leq \left[ \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i) \right]$$

$$\text{Loss} = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

$$\frac{1 = FP + BP}{=} \quad \boxed{5 \text{ epoch} = 5 \text{ times}}$$

Stochastic G.D.  $\Rightarrow$



SGD  $\Rightarrow$  Single Point  
N  $\Rightarrow$  [N iteration]

1 epoch  $\Rightarrow$  F.P + B.P  
Entire Data

5 epoch  $\rightarrow$  5 iteration  $\neq$  5 \* 5 = 25

Summary  $\rightarrow$  Vanilla GD  $\Rightarrow$  Entire Data  $\Rightarrow$  NN

1 epoch = F.P + B.P  
Entire Data

Single row

SGD  $\Rightarrow$  One Point at a time

1 epoch  $\Rightarrow$  n iteration

- n time weight will be updated

1

Which is faster?

Vanilla G.D.

2

Memory consumption

I.G.D  $\Rightarrow$  Vanilla G.D  $\uparrow \Rightarrow$  high memory  
RAM.  $\Rightarrow$  S.G.D.  $\downarrow$

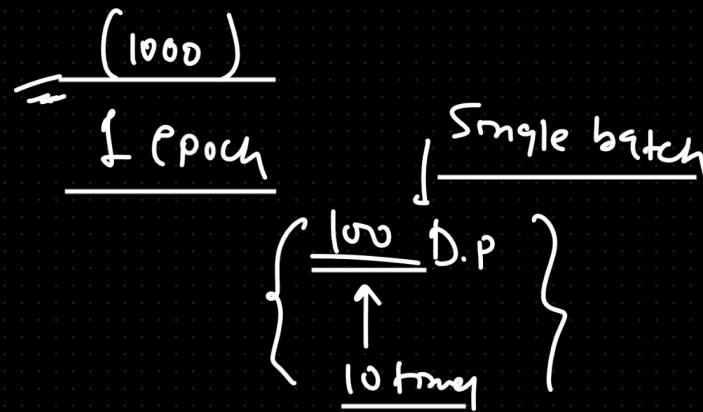
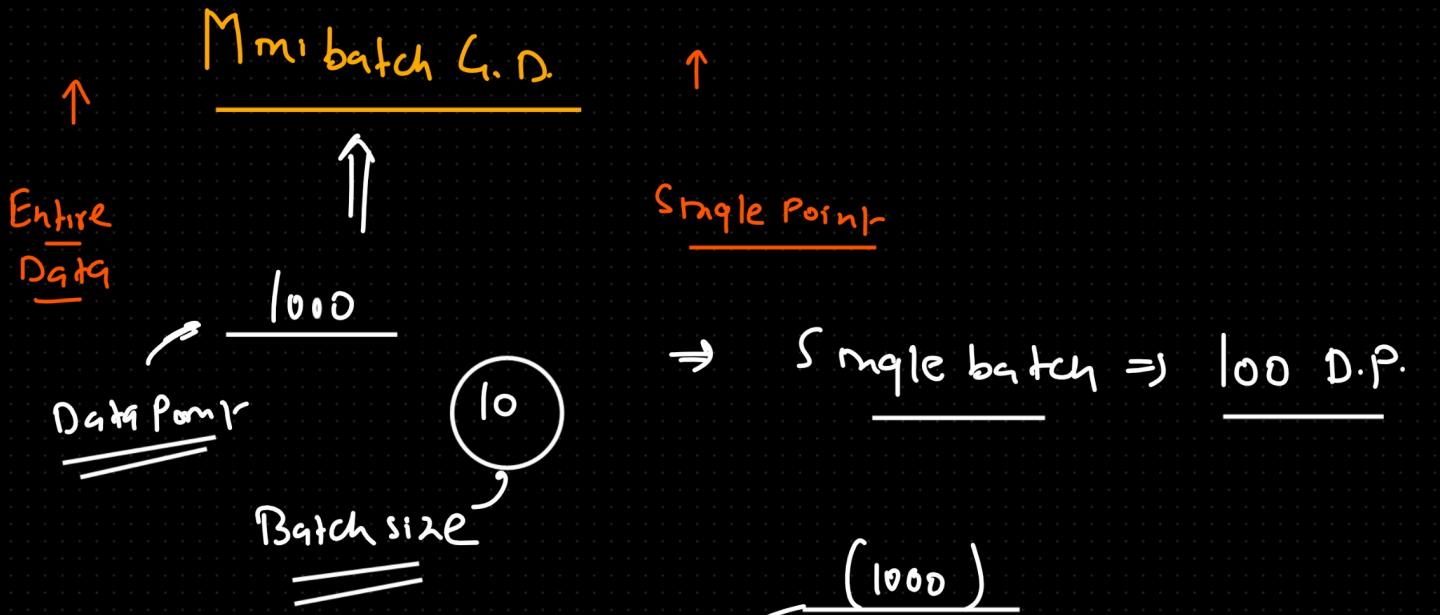
$\Rightarrow$  low memory

one point at a time

3

Good accuracy

S.G.D  $\Rightarrow$   $l_0 \times h \Rightarrow$   $n \times o$   $\Rightarrow$  weights  $\downarrow$  no. of weights  $\downarrow$  to update



GD, SGD, Minibatch G.D

- { ① epoch. }
- { ② Batch size }
- { ③ iteration }

SGD

1000

SGD

1000  $\Rightarrow$  NN

1 epoch  $\Rightarrow$  1000  $\xrightarrow{FP}$   
B.P.

5 epochs  $\Rightarrow$  5 times

epoch  
Batch size  
Iteration

Minibatch

~~entire~~

~~single~~

Batch

Batch

1000  $\approx$  10  $\Rightarrow$  100

(1000)  
= 2 epoch  $\Rightarrow$  100  
100  
100  
100  
100

(1000)  
= 2 epoch  $\Rightarrow$  100  
100  
100, 100, 100

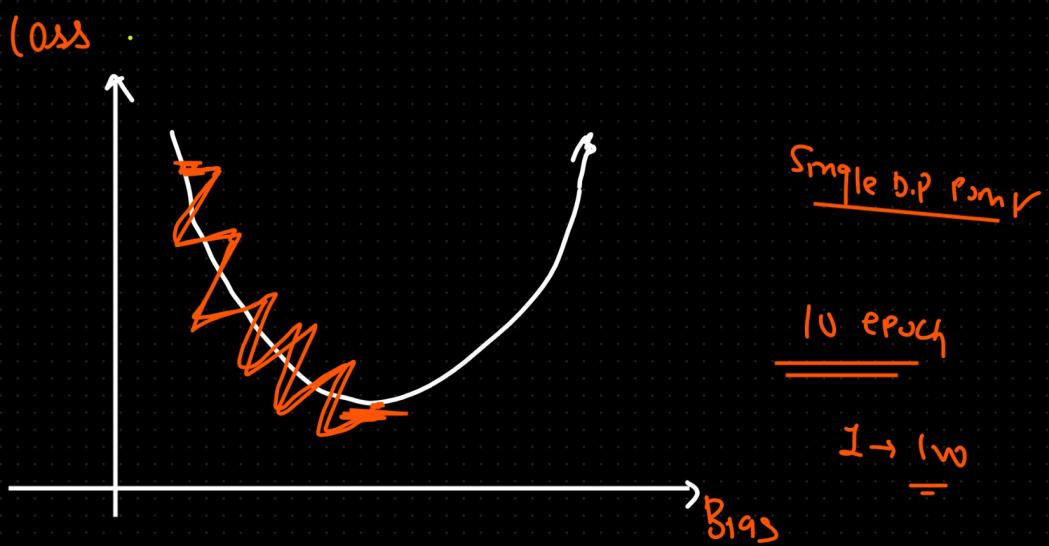
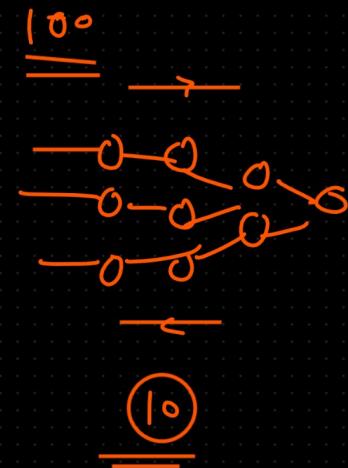
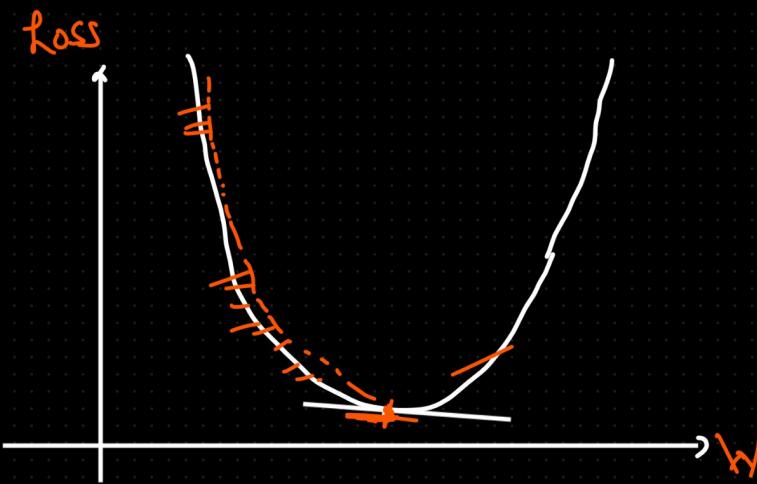
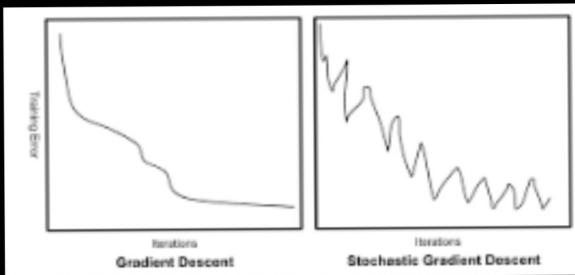
1000

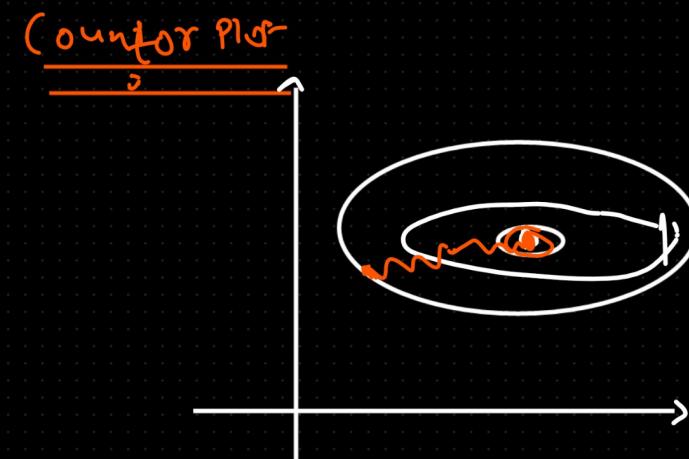
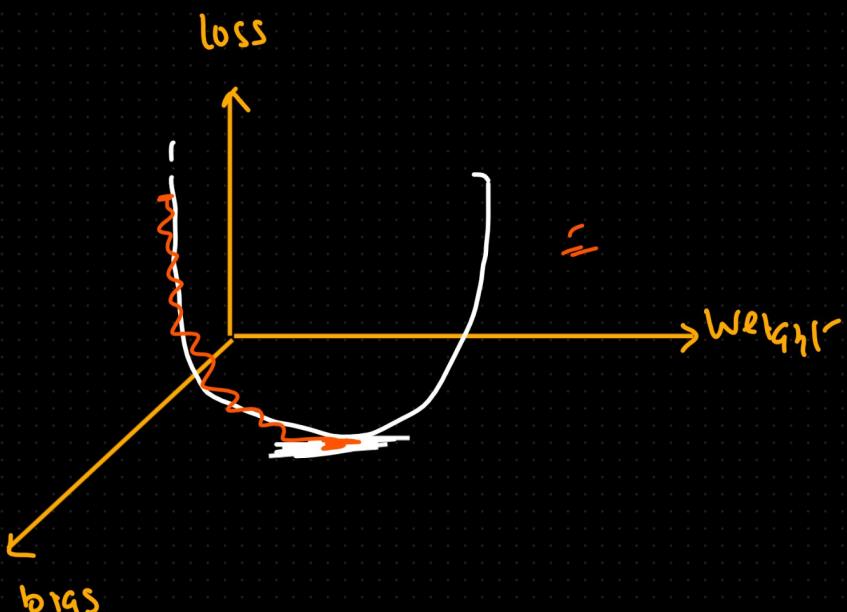
SGD

1 Data point

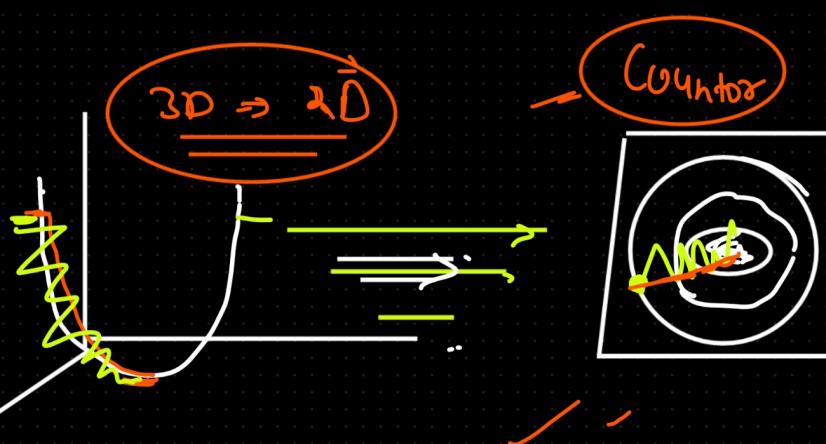
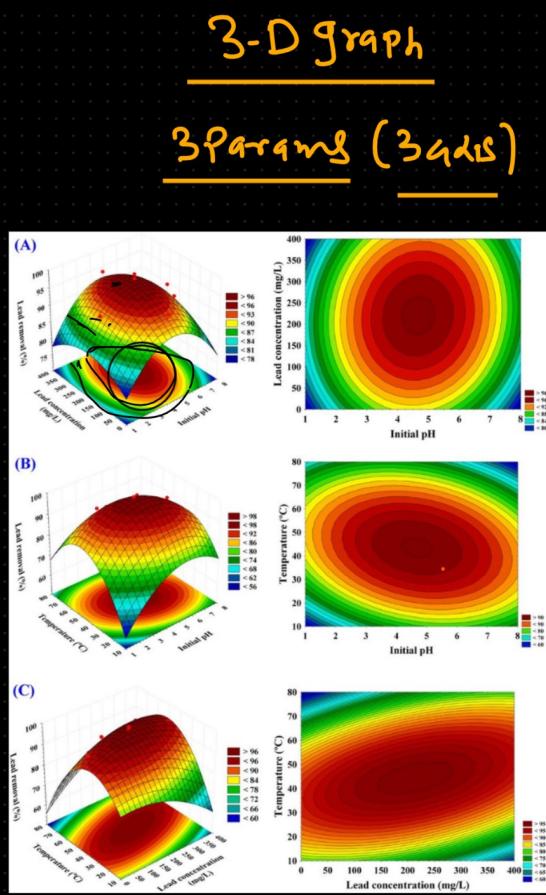
1 epoch  $\Rightarrow$  1

5 epoch  $\Rightarrow$  5  
1000  
5000





3-D Graph  $\Rightarrow$  2D

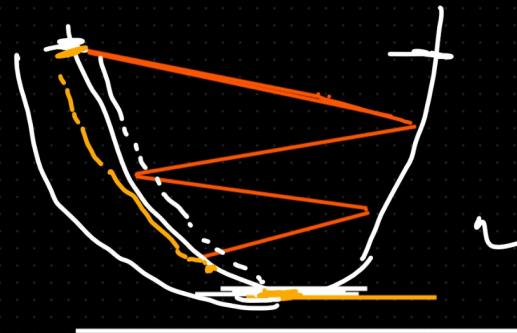


## Challenges with G.D.

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

$$\eta = 1, 2, 3, \dots$$

- 1 Learning rate
  - 2 Local minima Prob.
  - 3 Saddle Point
- $\equiv$

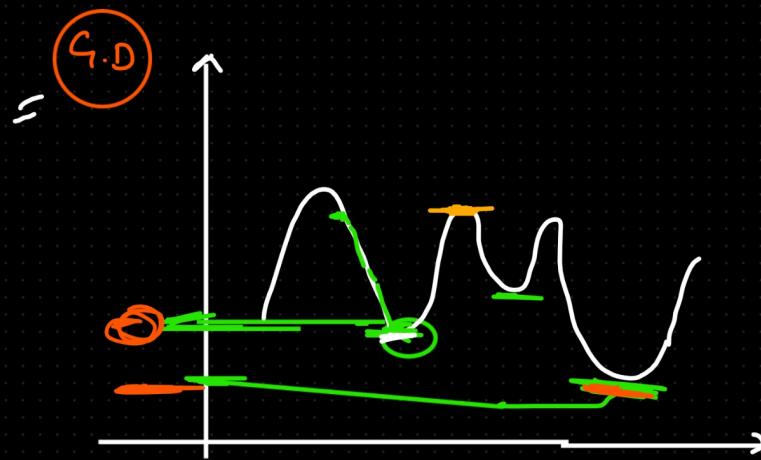


$$\eta = 0.001, 0.01$$

$$\eta = 0.1, 0.2$$

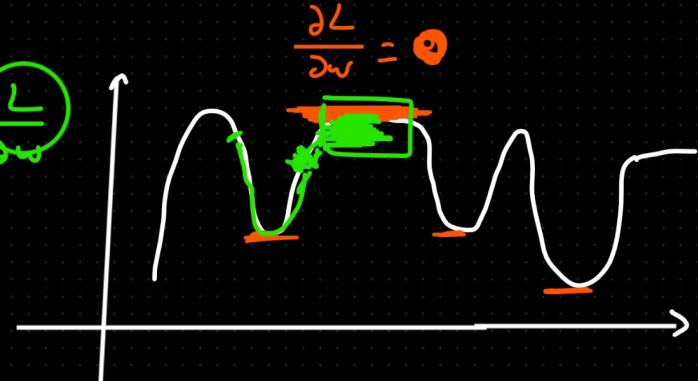
## Challenges G.D.D.

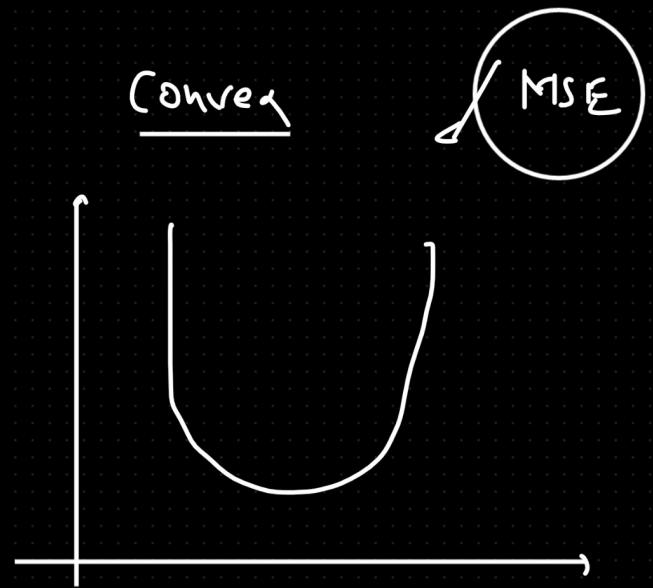
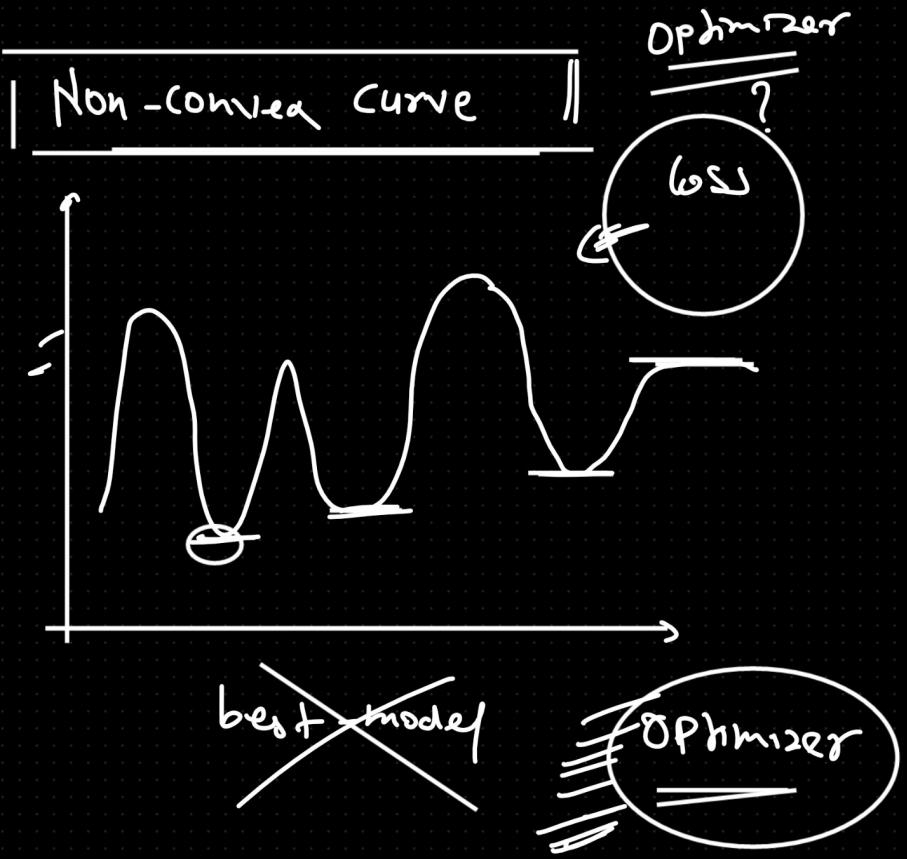
- 1 SGD with Momentum
- 2 Adagrad
- 3 NAG
- 4 RMS Prop
- 5 Adam



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

*No weight updation*





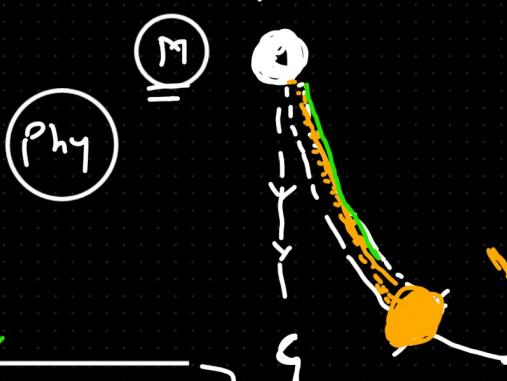
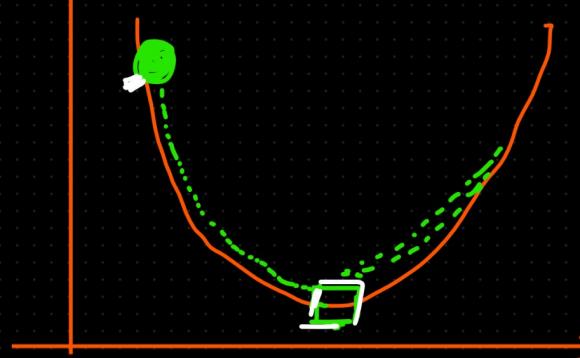
EWMA AND SMA

SGD with Momentum ?

Velocity

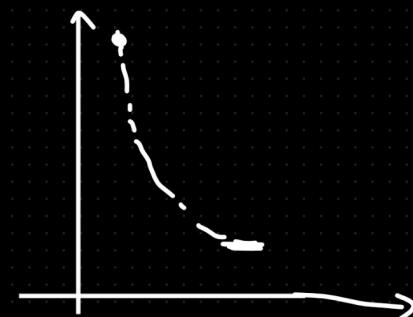
Max

Increase or Decrease



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

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



$$w_{t+1} = w_t - \gamma v_t$$

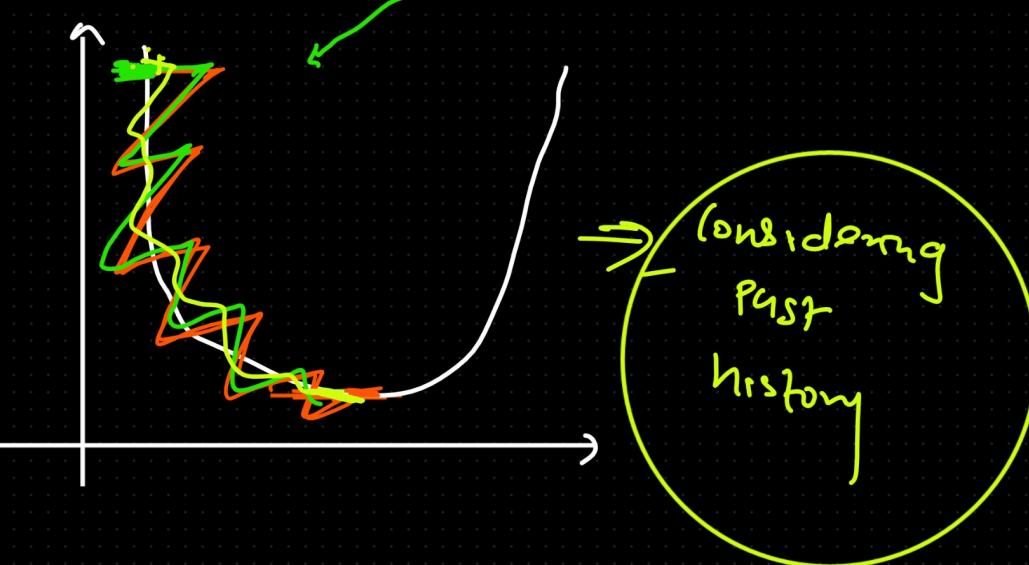
$$v_t = \beta \times v_{t-1} + \eta \frac{\partial L}{\partial w_t}$$

momentum

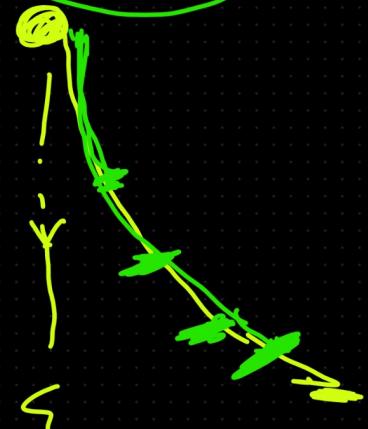
SGD      Mini-Batch

Noise

SGD with momentum



momentum



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

= Reducing the <sup>q</sup> Noise

Thursday → Keras tuner, tensorflow 2x

8-10

Sat, Sun → weights init-  
= batch norm, reg, dropout

(NN) → PyTorch

SMN



TS

⇒ Time dependent tuning