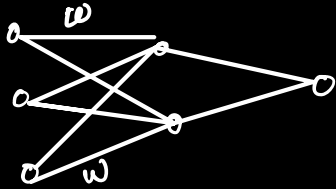


# Today's Agenda

1) Activation Functions

2) Optimizers



$$f = \boxed{wx + b} + \textcircled{\Gamma}$$

$$f(x) = x$$

Not Activating / Linear

Linear 5%

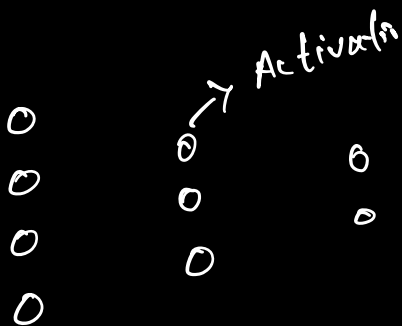
Non Linear 95%

Linear Regression

Activating  $\rightarrow$  Mathematical Transformation

$$\boxed{wx + b}$$

{ linear } Act fun



$\rightarrow$  Hyperparameters

1) Before the output in all layers

2) Output layer  $\rightarrow$  Rules are fixed

# Binary Classification

Rules { Last Activation  $\rightarrow$  Sigmoid  
Loss Function  $\rightarrow$  BCE

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

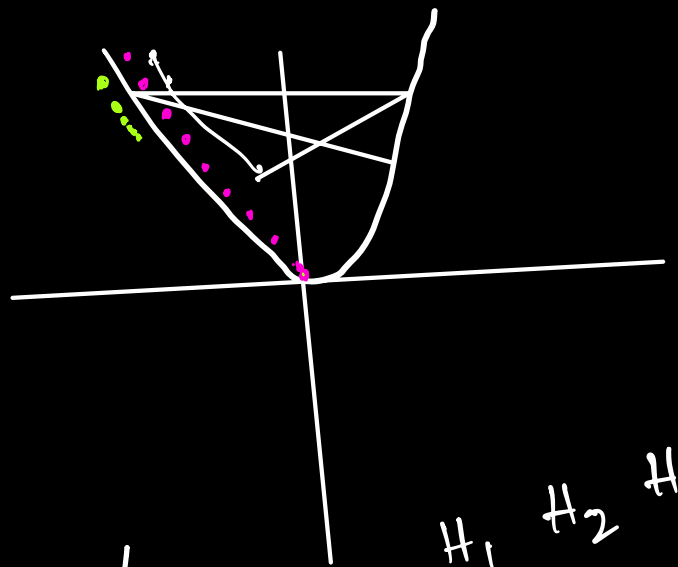
$$w_{\text{new}} = w_{\text{old}} - 0.00000065$$

$$w_{\text{new}} \approx w_{\text{old}}$$

No Weight Update  
No Training

Vanishing Gradient  
Problem

Loss will become  
stagnant far from  
minima.



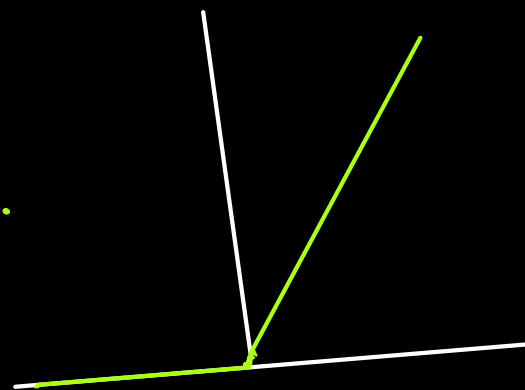
For a particular hidden layer

- ① All gradients are +ve.
- ② All gradients are -ve.

$H_1 \quad H_2 \quad H_3$   
 $H_1 \rightarrow +ve$   
 $\quad \quad -ve$   
 $0 \rightarrow$   
 $0$   
 $0$   
 $0$

## Relu Activation Fun<sup>c</sup>

✓ Negative Information  
is dumped  
+ve information is passed.  
Information loss



$$f(x) = \max(0, x)$$

$$z = wx + b$$

$$= z * \text{Act. Fun}^c$$

$$= z * 0$$

$$= 0$$

Dead Neuron

## Optimizers

1) Increase the performance of network

→ W/B update ~~\*~~

→ Act Fun<sup>c</sup> update

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

## Types of Optimizers

- 1) Gradient Descent → 1) Batch Gradient Descent  
 2) Stochastic Gradient  
 3) Mini Batch Gradient Descent

## Batch Gradient Descent

1000 → Data Points  
 Rows in Data

Epochs / Iterations

↓  
 complete dataset  
 seen by the  
 network

Epoch 1 {  $\frac{1000 \text{ datapoints}}{1 \text{ iteration}}$  }  
 weights will update

Epoch 2 {  $\frac{1000 \text{ datapoints}}{1 \text{ iteration}}$  }  
 weights will update

Epoch the big  
 cycle  
 Iteration  
 ↳ small  
 cycle

Pass all data points at once and then  
 weight updation

8GB RAM



1 billion data  
 points

\* OOM \*

Out of Memory  
 Issue

Advantages

Disadvantages

1) Convergence will happen

1) Resource Intensive  
Huge RAM

1 Epoch = 1 Iteration

CPU → RAM  
GPU → VRAM

Example :- 1000 datapoints  $\rightarrow$  100 Epochs

Weight Update  $\rightarrow$  100 times it will update

How many <sup>total</sup> weights  $\rightarrow \frac{1000 \times 100}{}$   
 ^  
 updation  $\Rightarrow 1000,00$

# Stochastic Gradient Descent

1000 data points

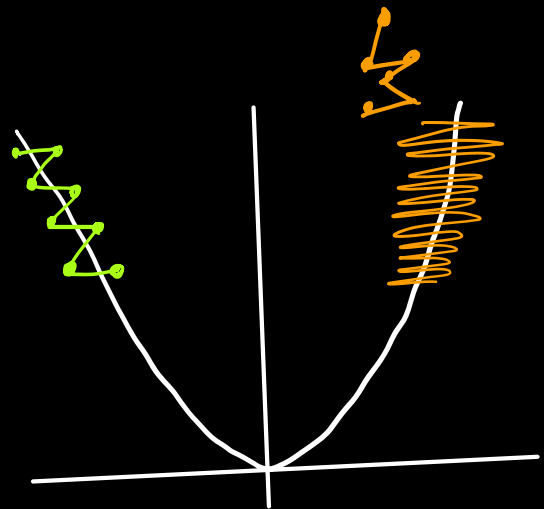
Epoch 1  $\left\{ \begin{array}{l} \frac{1 \text{ iteration}}{1 \text{ data point}} \\ \vdots \\ \frac{1000 \text{ iteration}}{1000 \text{ datapoint}} \end{array} \right.$

## Dis Advantages

i) Time taking

a) Resources Wastage

### 3) Jittery / Noise



## Advantage

1) No Resource Issue

## Mini Batch Gradient Descent

1000

BS  $\rightarrow$  90

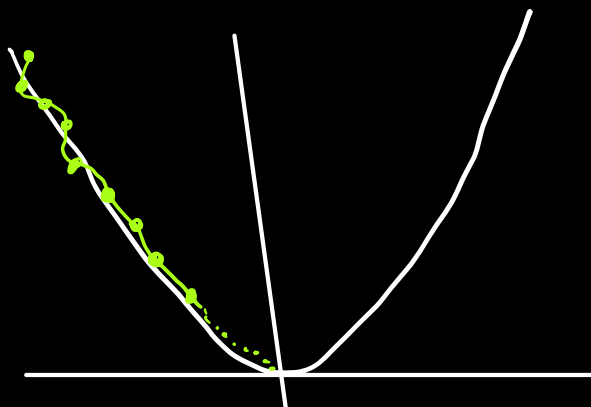
Iterations  $\rightarrow$

990 = 11 iteration

$\hookrightarrow$  10 datapoints

12<sup>th</sup> iteration  $\rightarrow$  10

Epoch 1 {  $\frac{\text{iteration 1}}{\text{Batch Size } 2}$



## Batch Size

1000 datapoints

BS  $\rightarrow$  100

Iterations  $\rightarrow$  10

$$= \frac{\text{Total data points}}{\text{Batch Size}}$$

if integer  
then  
iterations  
otherwise

int + 1

## Advantages

- 1) less Noise
- 2) Convergence will happen
- 3) Resource Efficient
- 4) Time less (sGD)

## Disadvantages

- 1) Noise exists

Batch Size  
↳ hyperparameter