

## Drop Out & regularization

\* Very <sup>deep</sup> huge artificial neural network have huge weight, huge bias then happens and tend to overfit the data

\* ~~also~~ underfit will happen only single layer neural network

\* underfit never happen in deep neural network

\* overfit never happen in single layer neural network

overfit  $\rightarrow$  high variance

① regularization  
 $L_1, L_2$

Random Forest

DT  $\rightarrow$  overfitting

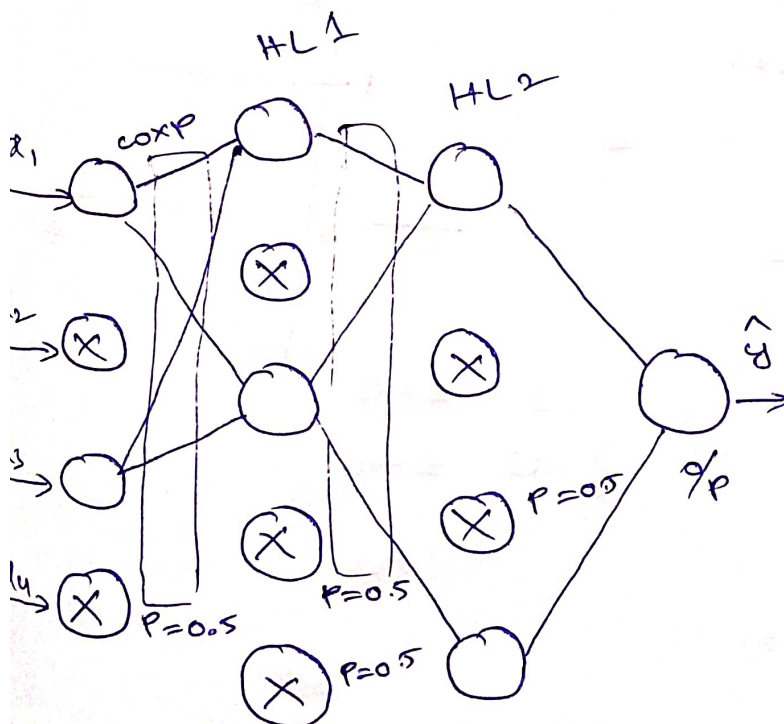
Subset of features

② Drop out

written

① Nitish Srivastava

② Geoffrey Hinton



Drop out ratio

$$0 \leq p \leq 1$$

\* select subset of features

\* randomly select neuron, activation function then deactivate it in forward and backward propagation

\* everytime select randomly

train data

test data  $\times p$

$p \rightarrow$  select by using hyper param tuning

\* Error is decreases

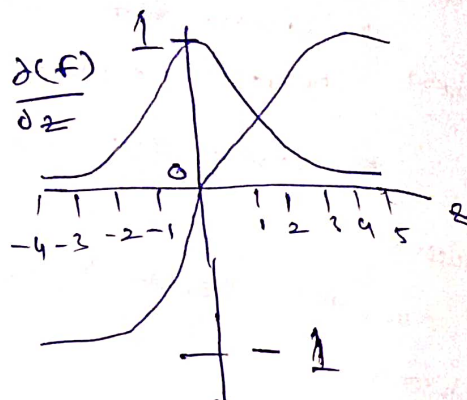
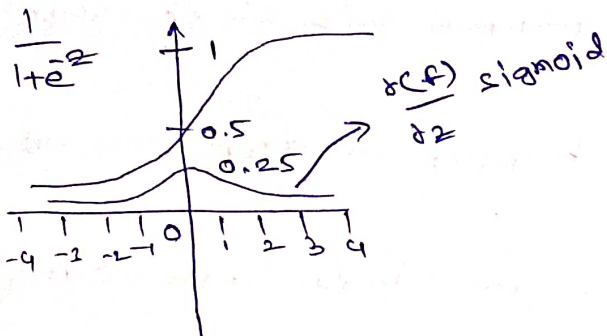
# Rectified linear Unit

$$0 \leq \frac{\partial(f)}{\partial(z)} \leq 0.25$$

$$\tanh(x) = \frac{1 - e^{-2x}}{1 + e^{-2x}}$$

① Sigmoid AF

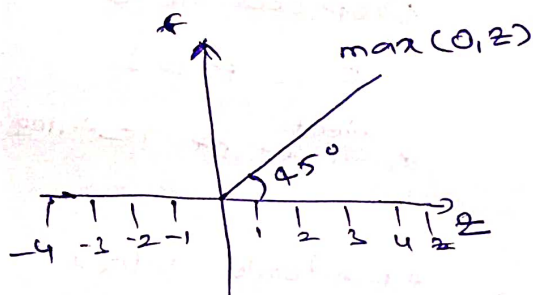
② Threshold Activation



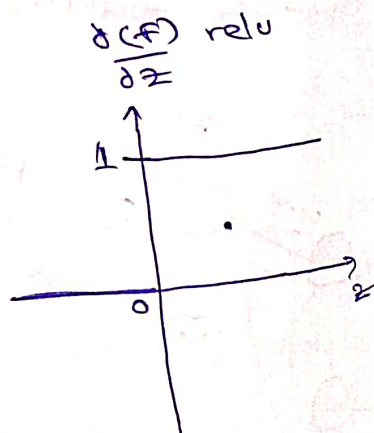
$$w_{new} = w_{old} - \eta \frac{\partial L}{\partial w}$$

$$0.25 \times 0.10 \times 0.001$$

relu :



$\Rightarrow$



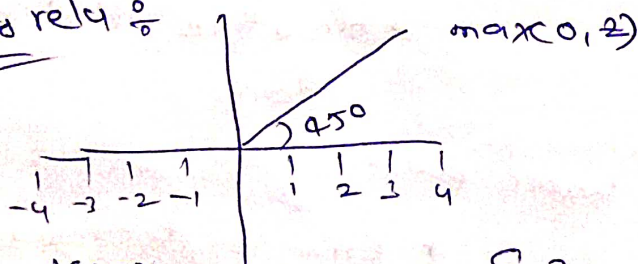
$$w_{old} = w_{new} - \eta \frac{\partial L}{\partial w}$$

1x1x1  
0x1x1

$\hookrightarrow$  dead neuron

$$\begin{cases} 1 & z > 0 \\ 0 & z < 0 \end{cases}$$

leaky relu :



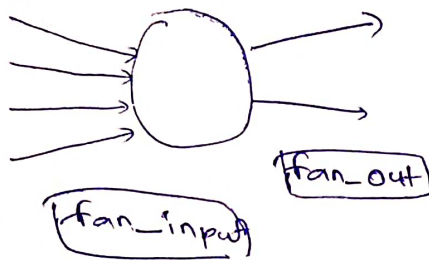
$$\frac{\partial(0.01)z}{\partial z} = 0.01$$

$$\begin{cases} z & z > 0 \\ 0.01(z) & z < 0 \end{cases}$$

# Weight Initialization

## Key Points

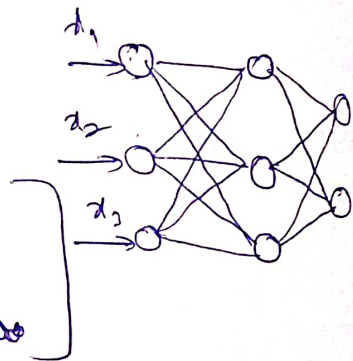
- ① weight should be small
- ② weights should not be same because each neuron will learn new
- ③ weights should have good variance



## Weight Initialization Techniques

- ① Uniform distribution

$$w_{ij} \sim \text{Uniform} \left[ \frac{-1}{\sqrt{\text{fan\_in}}}, \frac{1}{\sqrt{\text{fan\_out}}} \right]$$



- ② Xavier / Glorot

- ① Xavier Normal

$$w_{ij} \sim N(0, \sigma)$$

$$\sigma = \sqrt{\frac{2}{(\text{fan\_in} + \text{fan\_out})}}$$

- ② Xavier Uniform

$$w_{ij} \sim U \left[ \frac{-\sqrt{6}}{\sqrt{\text{fan\_in} + \text{fan\_out}}}, \frac{\sqrt{6}}{\sqrt{\text{fan\_in} + \text{fan\_out}}} \right]$$



① He init

① He uniform

② He Normal

$$w_{ij} \approx U \left[ -\sqrt{\frac{6}{fan_{in}}}, \sqrt{\frac{6}{fan_{out}}} \right]$$

$$w_{ij} \approx N(0, \sigma)$$

$$\sigma = \sqrt{\frac{2}{fan_{in}}}$$

## Stochastic Gradient Descent

1000 data points

$$w_{new} = w_{old} - \eta \times \left[ \frac{\partial L}{\partial w_{old}} \right] \left( \frac{\partial L}{\partial w_{old}} \right) \rightarrow n \text{ data point} \rightarrow \text{Gradient Descent}$$

↳ 1 data point → SGD

↳ 10 data point → mini batch SGD  
100 (sample)

$$\text{mini batch SGD} = \text{loss} = \sum_{i=1}^K (y_i - \hat{y}_i)^2$$

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

$$\text{SGD} = \text{loss} = (y - \hat{y})^2$$

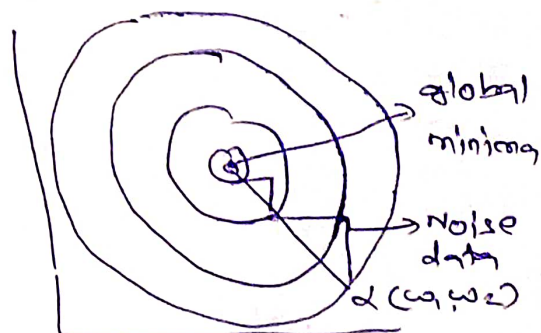
$$w_{new} = w_{old} - \eta \times \frac{\partial L}{\partial w_{old}}$$

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

mini batch  
SGD

↳ population

↳ sample



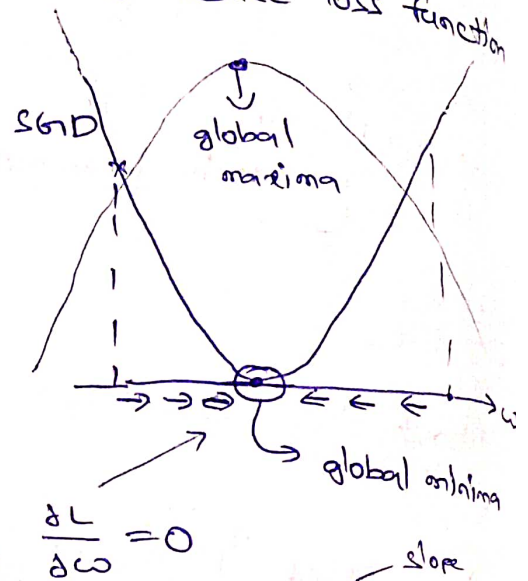
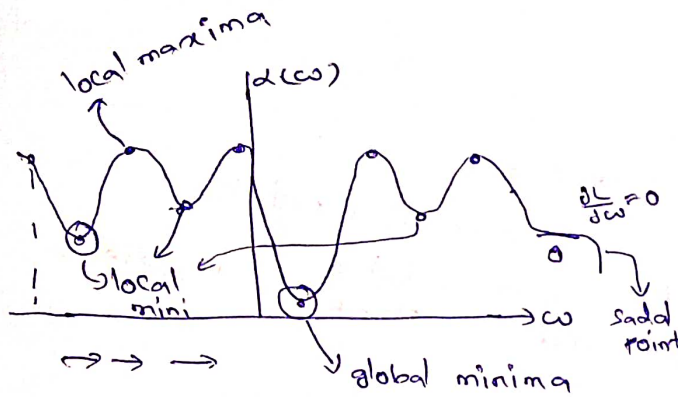
\* remove noise stochastic GD with momentum

# Global minima Vs local minima.

① GD, SGD, Batch SGD  $\Rightarrow$  optimizers used

for reduce loss function

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

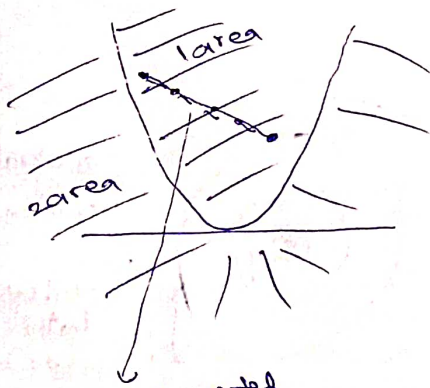


$$w_{new} = w_{old} - \eta \left( \frac{\partial L}{\partial w_{old}} \right)$$

slope  
 $\uparrow$

## ① Convex function :

linear regression  
logistic

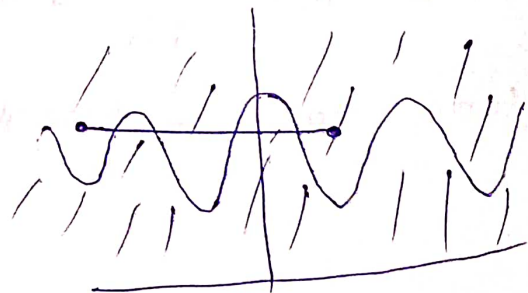


if point present in same curve is called Convex fun

\* occur in ml problem  
tangent

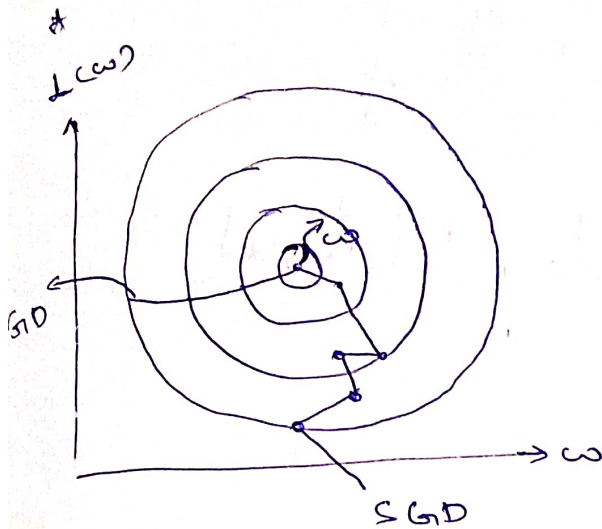
## ② Non Convex function :

↳ deep learning



\* mostly occur in deep learning  
\* we choose two points and connect both. if both are present in same curve is convex other wise non convex

# Stochastic Gradient Descent with Momentum:



## GD optimizer:

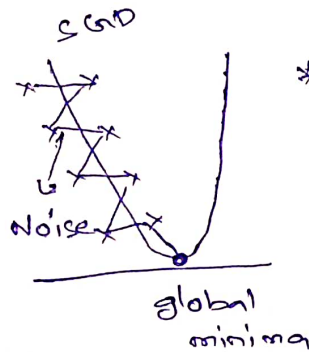
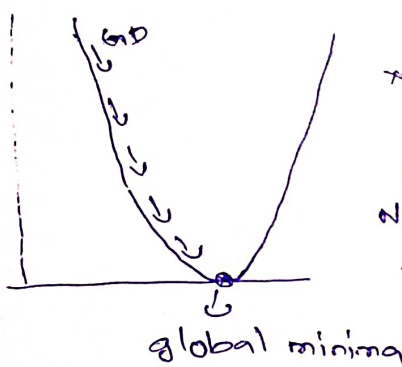
\* I am take all the datapoints and it will converge fastly

\* if data point is more it will take more time and Computational power overcome we use

SGD

## SGD:

\* select batch (like, 100, 200 data points select)



\* Noise data

\* remove noise using exponential moving average

time interval		<u>Exponent moving average</u>			
$t_1$	$t_2$	$t_3$	$t_4$	$t_5$	$\dots t_n$
$b_1$	$b_2$	$b_3$	$b_4$	$b_5$	$\dots b_n$

$$0 \leq \gamma \leq 1$$

$$0.5$$

$$V_{t_1} = b_1$$

$$V_{t_2} = \gamma \times V_1 + b_2$$

$$= 0.5 \times b_1 + b_2$$

$$V_{t_3} = \gamma \times V_2 + b_3$$

$$= \gamma \times (\gamma \times V_1 + b_2) + b_3$$

$$= \gamma^2 \times b_1 + \gamma b_2 + b_3$$

$$= 0.25 \times b_1 + 0.5 b_2 + 1 \times b_3$$