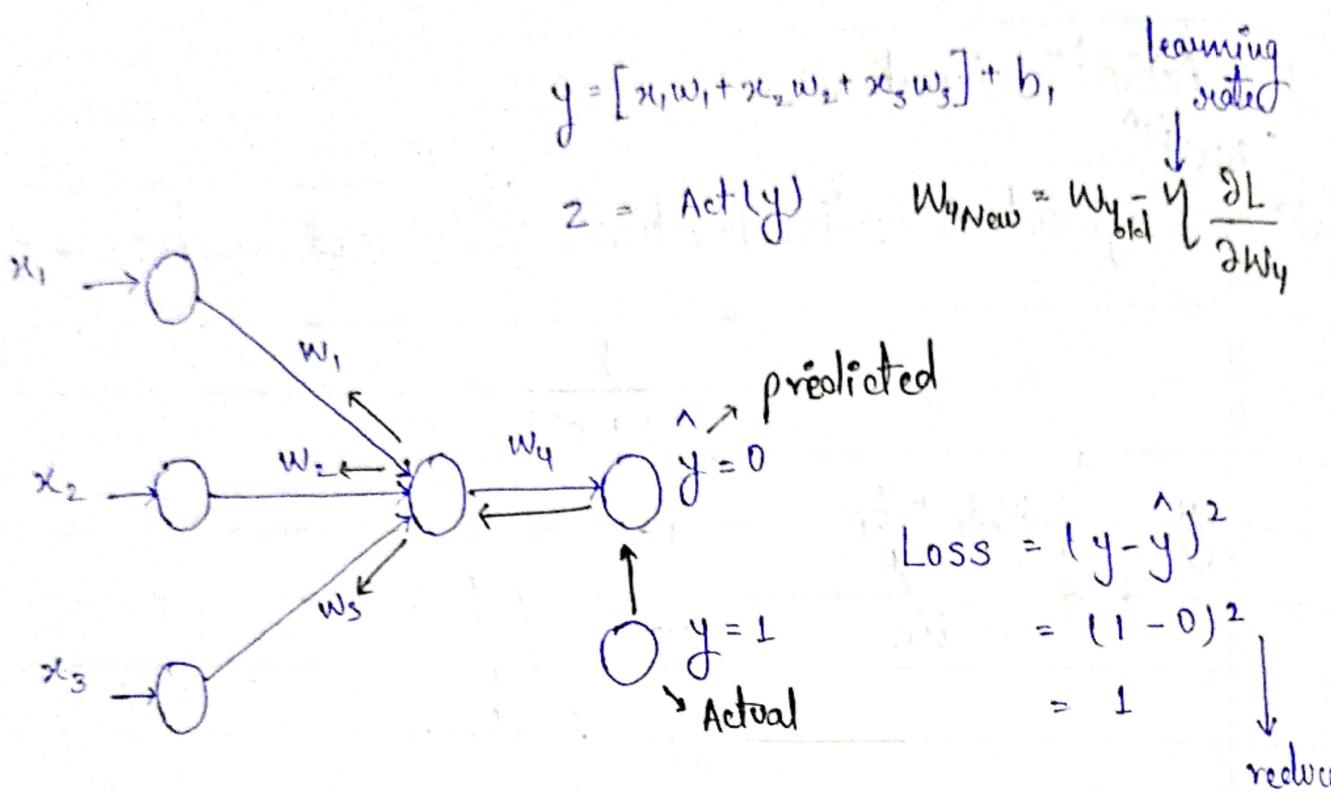


How to train Neural Network with Backpropagation



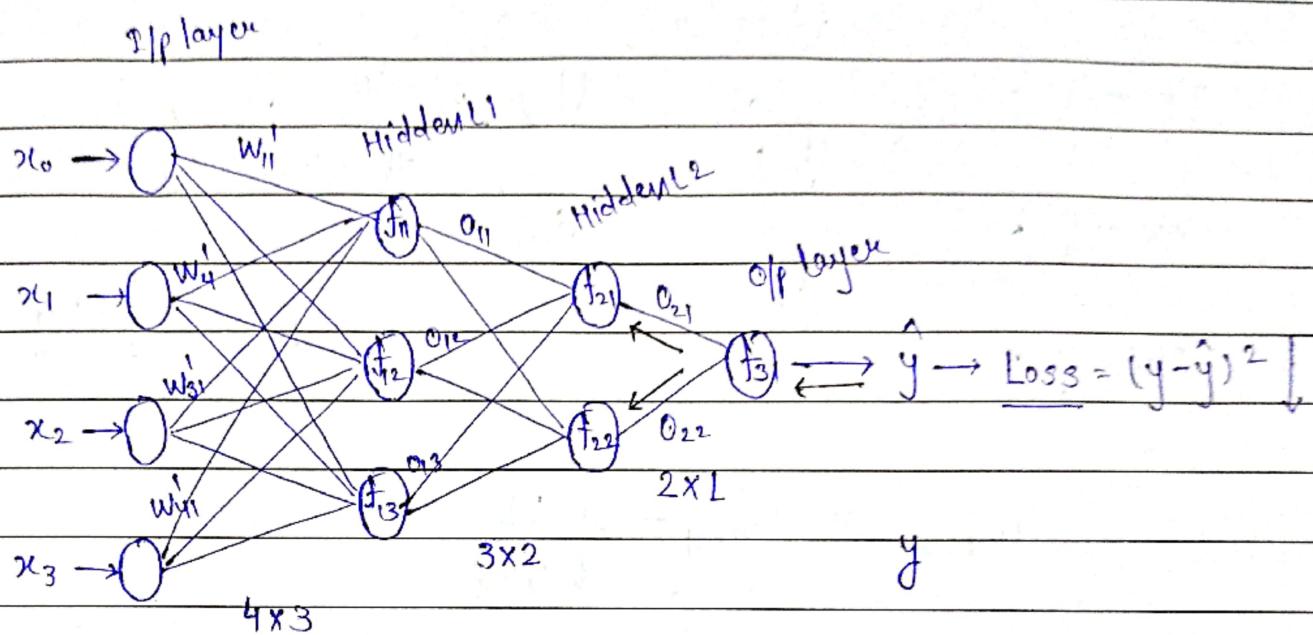
- We have to update weights during training phase to decrease loss function.
- ∴ To minimize loss, Optimizers are used.
- To update the weights Backpropagation is done.

Similarly,

$$w_{3,\text{New}} = w_{3,\text{Old}} - \eta \frac{\partial \text{Loss}}{\partial w_3}$$

learning rate

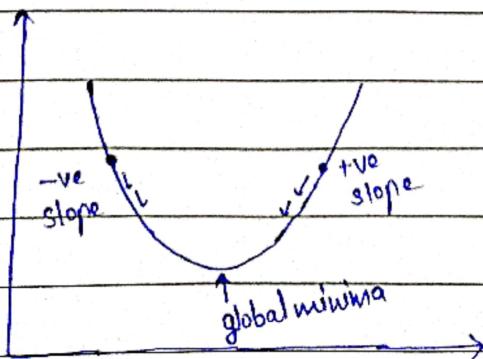
How to train Multilayer Neural Network & Gradient Descent:



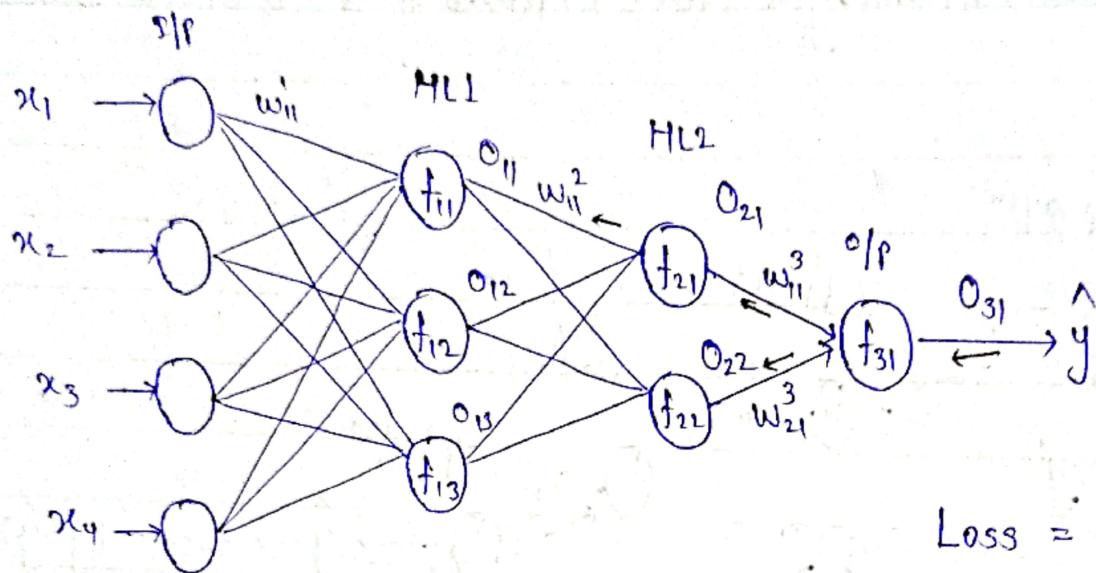
→ Weights Updation : $W_{\text{new}} = W_{\text{old}} - \eta \frac{\partial L}{\partial W}$

learning rate

Gradient Descent :



Chain Rule of Differentiation with Back Propagation:



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

$$w_{11}^3_{\text{new}} = w_{11\text{old}}^3 - \eta \frac{\partial L}{\partial w_{11}^3}$$

$$\frac{\partial L}{\partial w_{11}^3} = \frac{\partial L}{\partial o_{31}} \cdot \frac{\partial o_{31}}{\partial w_{11}^3}$$

$$\frac{\partial L}{\partial w_{21}^3} = \frac{\partial L}{\partial o_{31}} \cdot \frac{\partial o_{31}}{\partial w_{21}^3}$$

$$\frac{\partial L}{\partial w_{11}^2} = \frac{\partial L}{\partial o_{31}} \cdot \frac{\partial o_{31}}{\partial o_{21}} \cdot \frac{\partial o_{21}}{\partial w_{11}^2}$$

Vanishing Gradient Problem:

- The ~~difference~~, difficulty of training deep neural networks because the gradients, used to update the network's weights, become extremely small as they are propagated backward from the output layer to the earlier layers.

Root Cause

It typically arises when using sigmoid or tanh activation functions, which squash input into a small range, causing derivatives to be less than 1, and their repeated multiplication causes gradient to shrink exponentially.

- Solutions include using ReLU, proper weight initialization, batch normalization, and residual connections.

Exploding Gradient Problem:

- Gradients grow exponentially large during backpropagation in deep neural networks

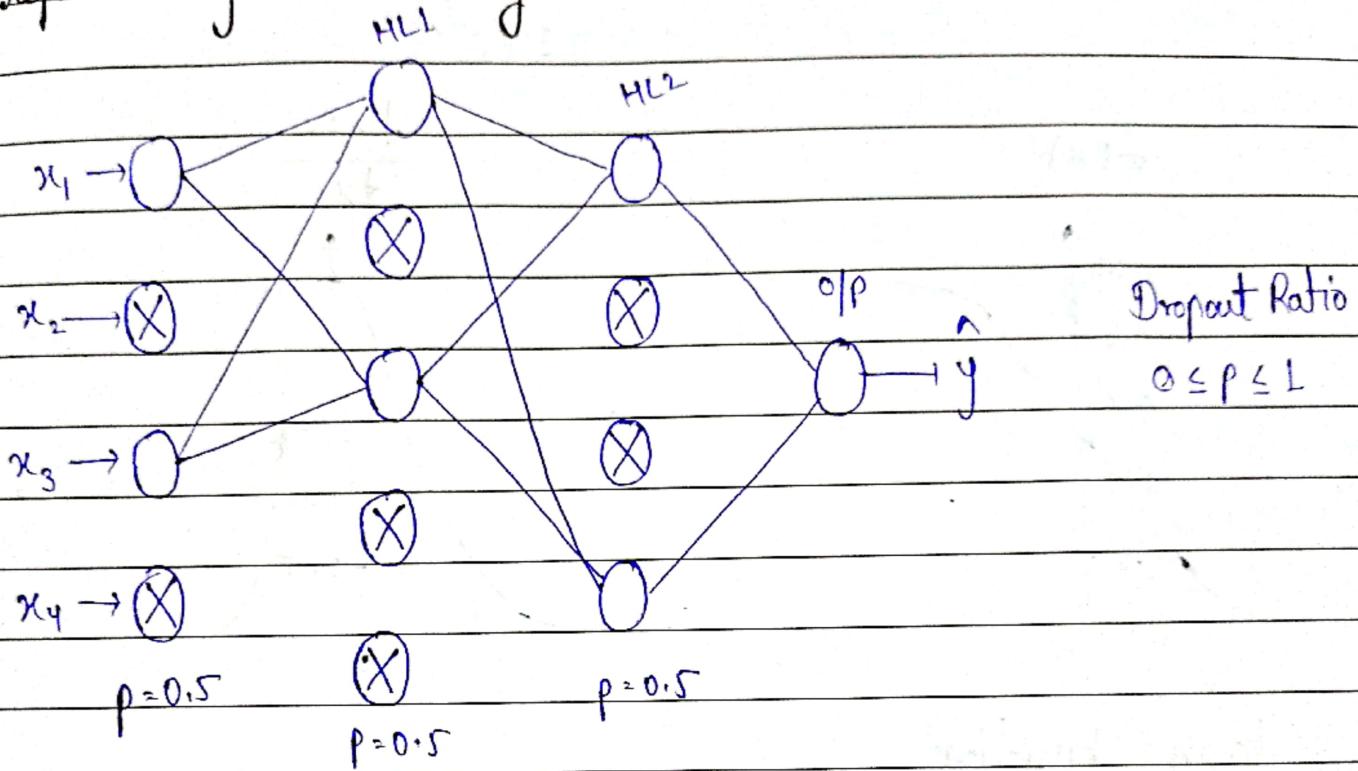
Root Cause

→ Caused by repeatedly multiplying large weights, especially in deep networks or RNNs.

→ The exploding gradient problem leads to unstable training by causing excessively large weight updates, often resulting in divergence or NaN values.

By applying techniques like Gradient clipping, proper weight initialization and normalization, this issue can be effectively managed, ensuring stable and successful model training.

Dropout layers and Regularization:



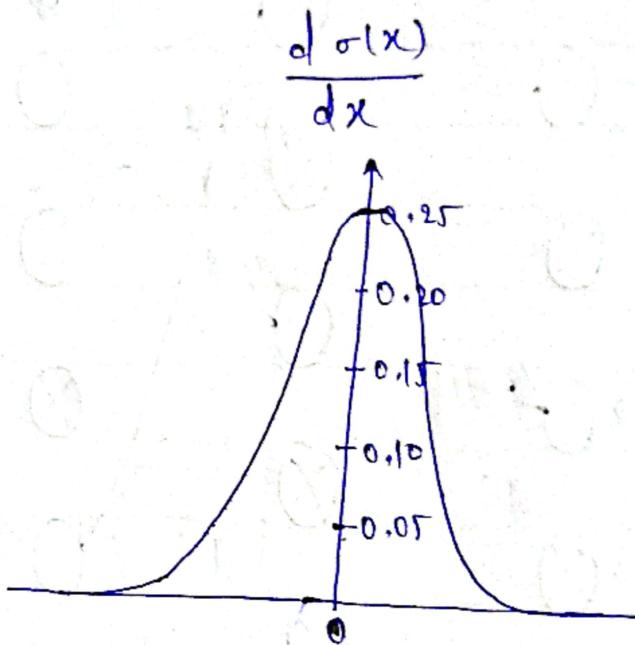
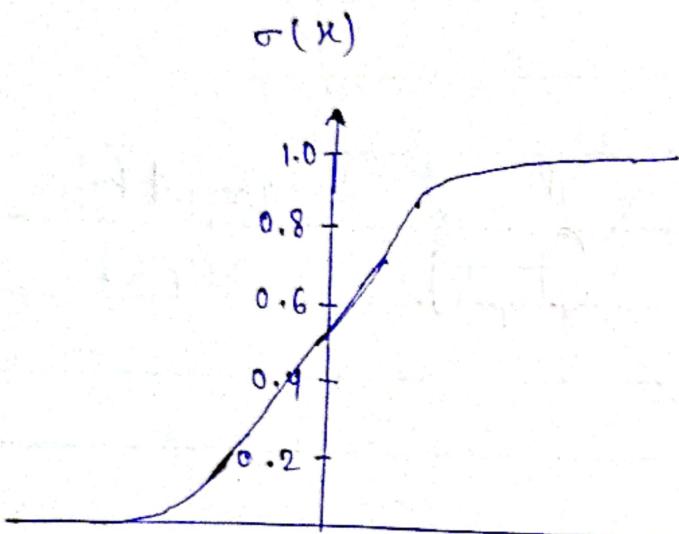
→ Dropout is a regularization technique used in neural networks to prevent overfitting.

- Randomly drop out a fraction of the neurons during training.
- During each training step, selected neurons are temporarily ignored.
- At test time, all neurons are used, but their outputs are scaled to account for the dropout.

Activation functions:

1. Sigmoid function:

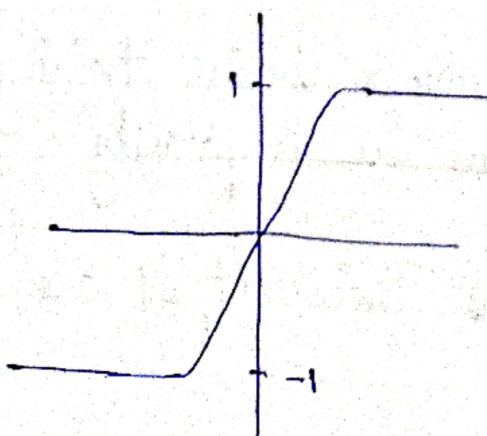
$$\sigma(x) = \frac{1}{1+e^{-x}}$$



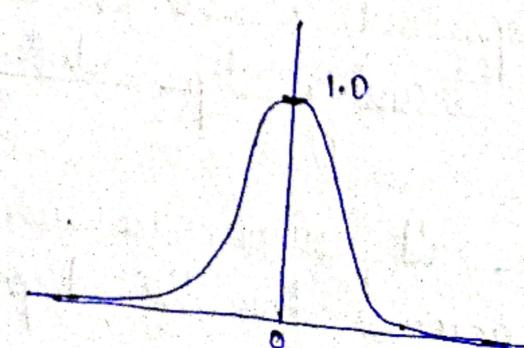
2. tanh function:

$$\tanh x = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$\tanh(x)$



$$\frac{d\tanh(x)}{dx}$$

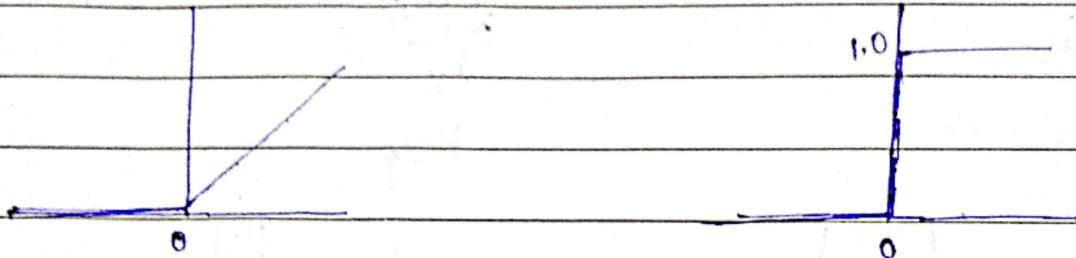


3. ReLU Function:

$$\text{ReLU} = \max(0, x)$$

$\text{ReLU}(x)$

$$\frac{d\text{ReLU}(x)}{dx}$$

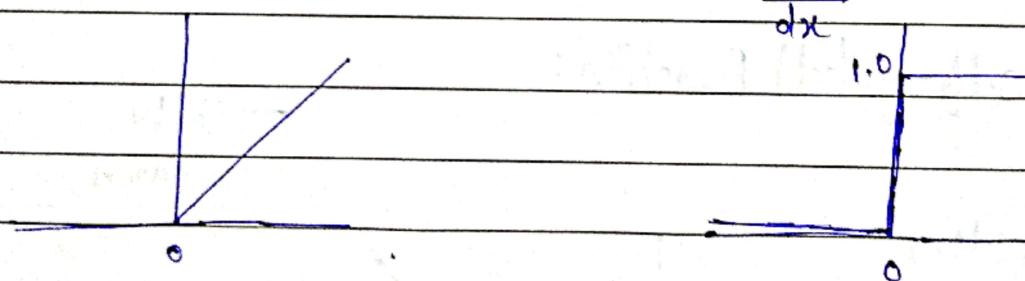


4. Leaky ReLU function:

$$f(x) = \max(0.01x, x)$$

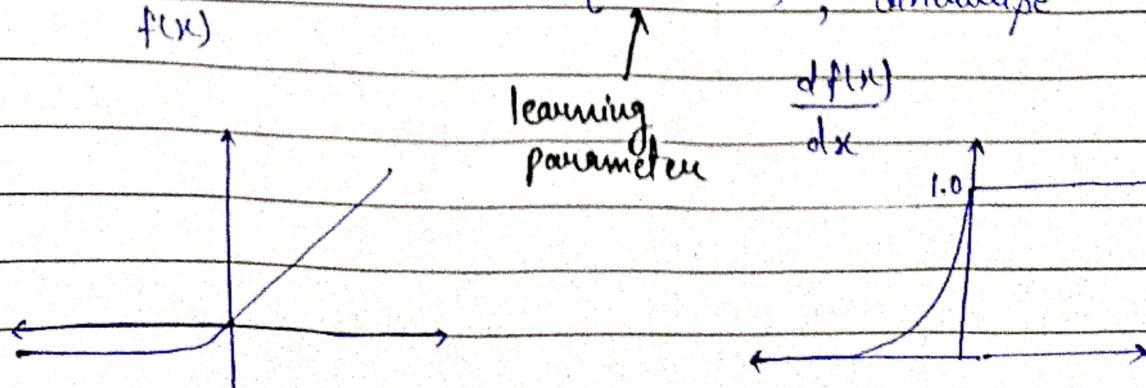
$f(x)$

$$\frac{df(x)}{dx}$$



5. ELU (Exponential Linear Units) Function:

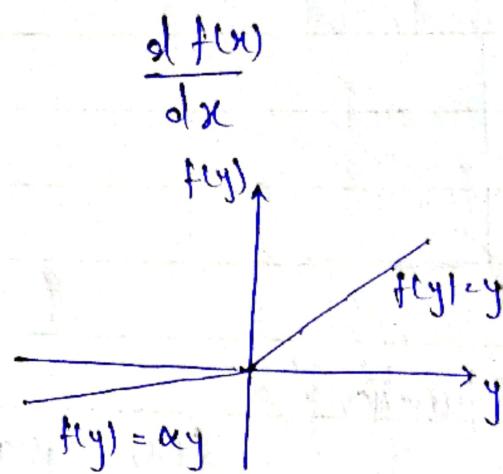
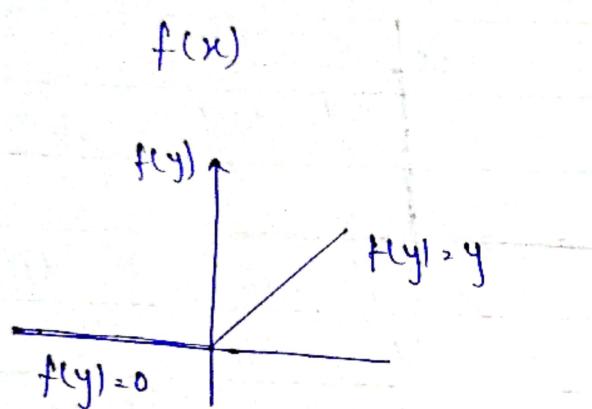
$$f(x) = \begin{cases} x, & x > 0 \\ \alpha(e^x - 1), & \text{otherwise} \end{cases}$$



6. PReLU (Parametric ReLU):

$$f(x) = \begin{cases} x & , x > 0 \\ \alpha * x & , x \leq 0 \end{cases}$$

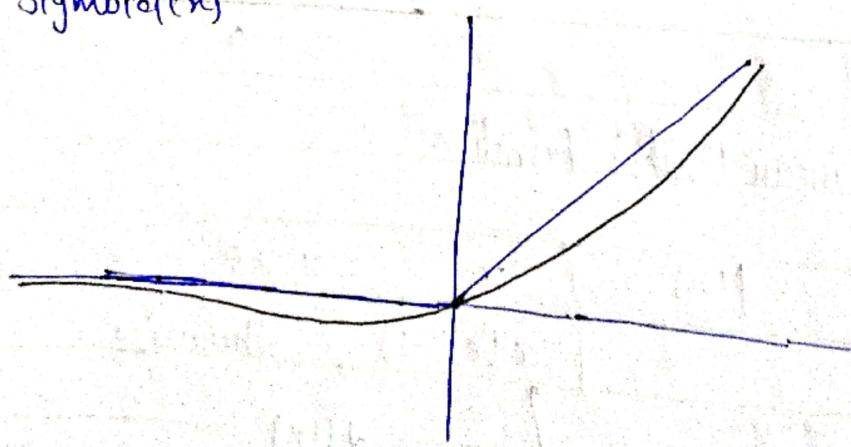
↑
learning parameter



7. Swish (A self-gated) Function:

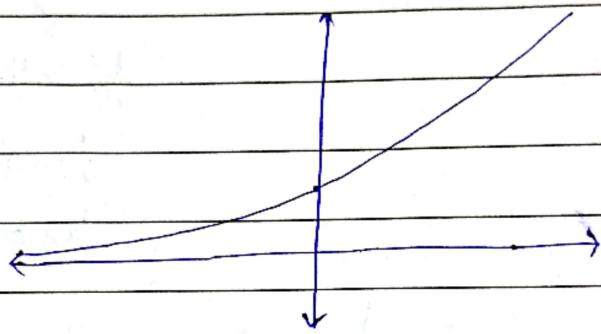
$$y = x * \text{Sigmoid}(x)$$

— ReLU
— Swish



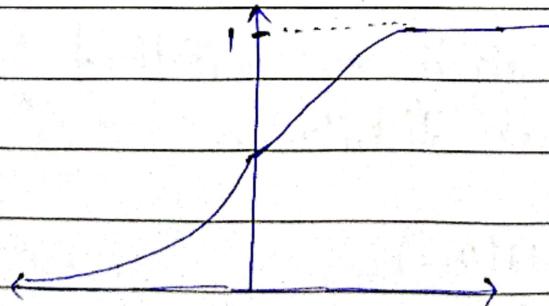
8. Softplus Function:

$$f(x) = \ln(1 + e^{x_1})$$



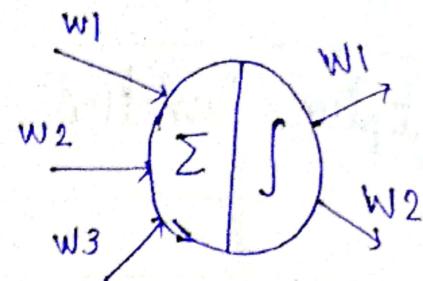
9. Softmax function:

$$\frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}}$$



Weight Initialization in Neural Networks:

- Weights should be small
- Weights should not be same
- Weights should have good variance.



for above neuron,

$$\text{fan-in} = 3$$

$$\text{fan-out} = 2$$

All the weights are assigned zero as the initial value. This is zero initialization. This kind of initialization is highly ineffective as neurons learn the same feature during each iteration.

2. Random Initialization:

a) Random Normal: The weights are initialized from values in a normal distribution.

$$w_i \sim N(0,1)$$

b) Random Uniform: The weights are initialized from values in a uniform distribution.

$$w_i \sim N(0,1)$$

3. Xavier/Glorot Initialization:

The weights are assigned from values of a uniform distribution as follows:

$$w_i \sim U \left[-\frac{\sigma}{\sqrt{\text{fan-in} + \text{fan-out}}}, \frac{\sigma}{\sqrt{\text{fan-in} + \text{fan-out}}} \right]$$

→ Xavier/Glorot Initialization often termed as Xavier Uniform Initialization, is suitable for layers where activation function used is Sigmoid.

4. Normalized Xavier/Glorot Initialization:

$$w_i \sim N(0, \sigma)$$

$$\sigma = \sqrt{\frac{G}{\text{fan-in} + \text{fan-out}}}$$

5. He Uniform Initialization:

$$w_i \sim U \left[-\frac{\sqrt{G}}{\sqrt{\text{fan-in}}}, \frac{\sqrt{G}}{\sqrt{\text{fan-out}}} \right]$$

→ He Uniform Initialization is suitable for layers where ReLU is used.

G. He Normal Initialization:

$$W_i \sim N[0, \sigma]$$

$$\sigma = \sqrt{\frac{2}{fan-in}}$$

→ Suitable when ReLU is used.

$$\frac{\partial L}{\partial W_{old}} \rightarrow N \text{ data points} \rightarrow \text{Gradient Descent}$$

Stochastic Gradient Descent

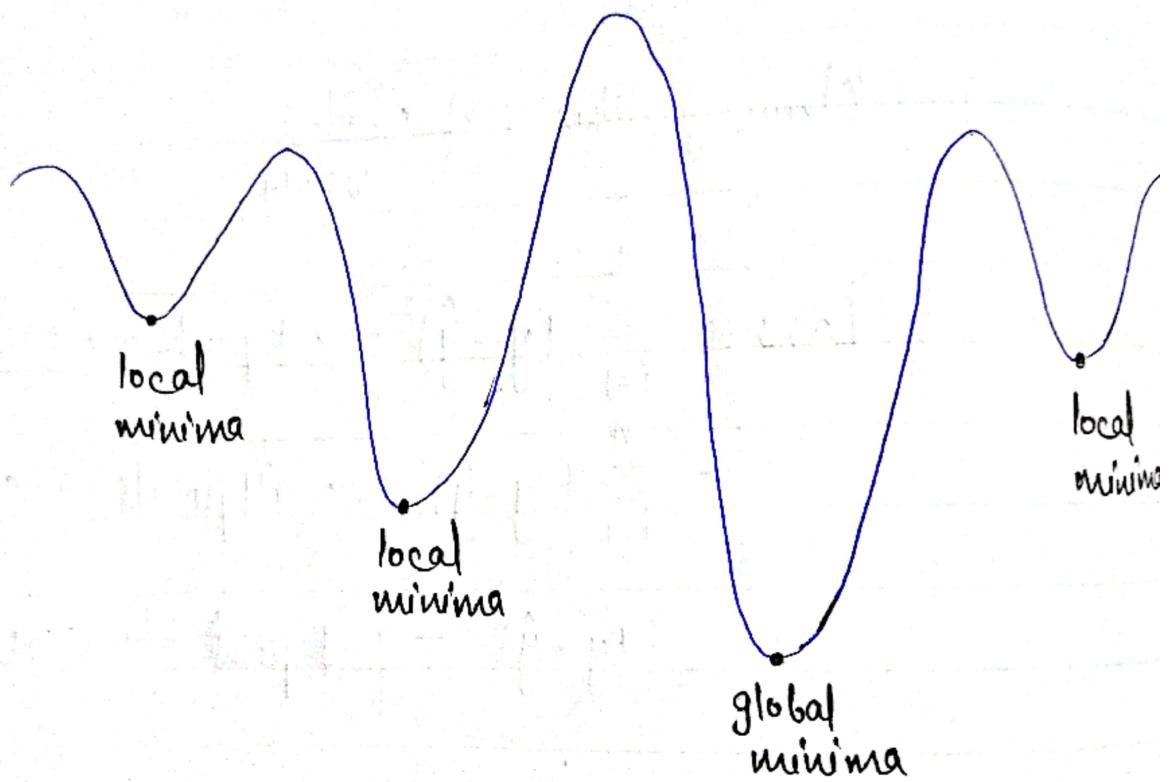
$$W_{new} = W_{old} - \eta \times \frac{\partial L}{\partial W_{old}}$$

$$\text{Loss} = \sum_{i=1}^k (y - \hat{y})^2 \rightarrow k \text{ points} \rightarrow \text{Mini Batch SGD}$$

$$= \sum_{i=1}^n (y - \hat{y})^2 \rightarrow \text{all points} \rightarrow \text{GD}$$

$$= (y - \hat{y})^2 \rightarrow 1 \text{ point} \rightarrow \text{SGD}$$

Global Minima vs Local Minima



- A global minimum is the absolute lowest point on the loss surface.
- A local minimum is a point where the loss is lower than neighboring points but not the lowest overall.