

Perceptron

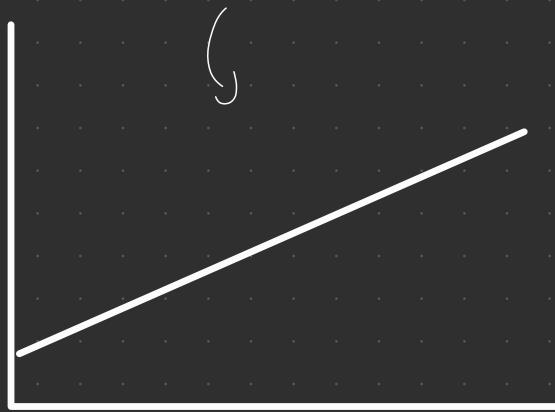
The Perceptron is the simplest type of artificial neuron, introduced by Frank Rosenblatt in 1958. It is the foundation of neural network and Deep learning.

It mimics how a biological neuron works taking inputs, processing them, and giving an output.

Table →	Age	cholesterol	B.P	heart Disease	output
	(28)	(150)	(110)		1
	36	120	90	0	
	42	180	160	1	

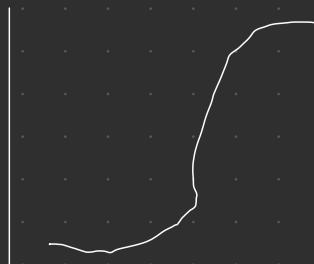
$$\text{Step -1} \quad y = mx + b \quad \rightarrow \quad h(x) = \theta_0 + \theta_1 x_1$$

$$h(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 \rightarrow$$



Step -2 → Sigmoid Activation function

$$\left(\begin{array}{l} \\ \end{array} \right) = \frac{1}{1+e^{-z}} \Rightarrow z = (\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3)$$



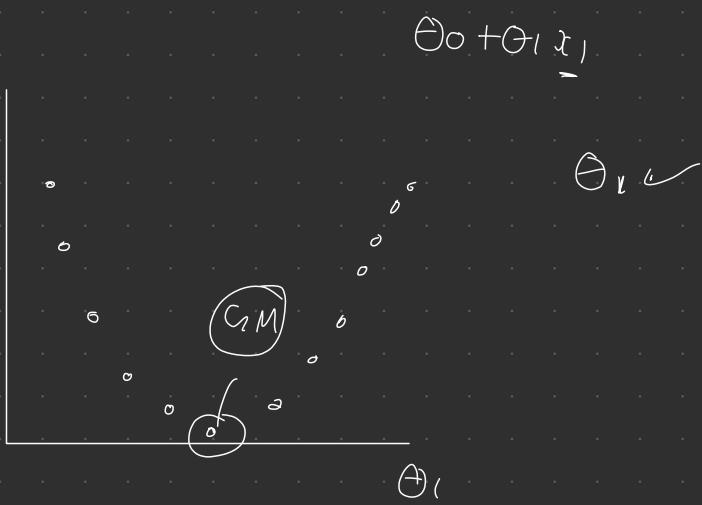
Step - 3 log loss function

\Rightarrow loss func

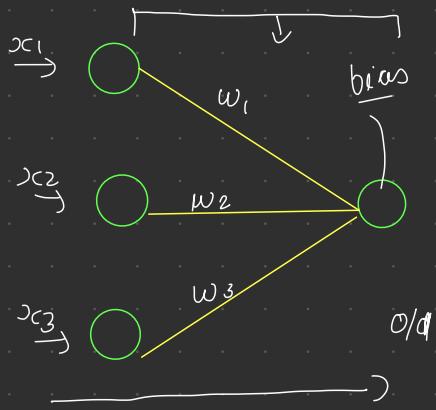
Step - 4 - Optimizer

$$\boxed{\theta_0 = 0}$$

log loss



We did all of this using perceptron



$$\begin{aligned} & \Rightarrow (x_1 \times w_1 + x_2 \times w_2 + x_3 \times w_3 + \text{bias}) \\ & \Rightarrow (\theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_0) \end{aligned}$$

same Sigmoid Activation function

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

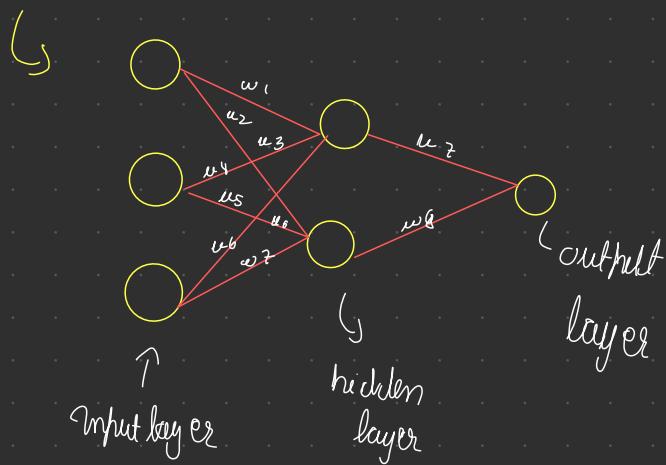
$$O/P \Rightarrow = O/1$$

Limitations of Perceptrons

\rightarrow Step Activation function

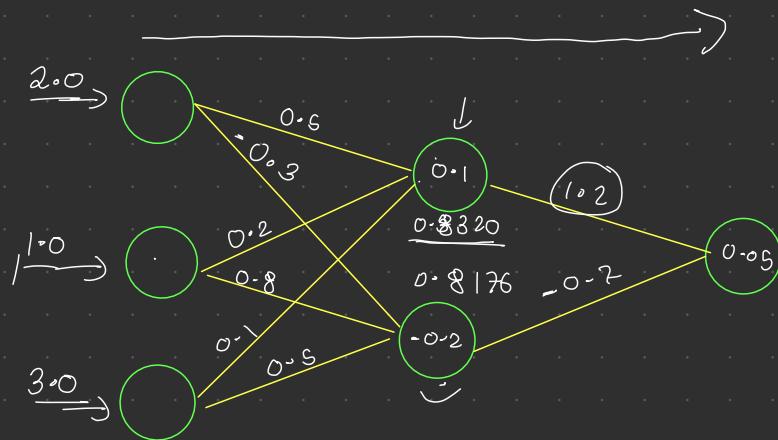
$$\boxed{O/1}$$

ANN (Artificial Neural Network)



Forward Propagation

$$\begin{array}{c|c|c|c} x_1 & x_2 & x_3 & y \\ \hline 2.0 & 1.0 & 3.0 & 1 \end{array}$$



$$\begin{aligned} h_{\theta_1} &= (x_1 w_1) + (x_2 (w_2) + (x_3 (w_3)) + b \\ &= (2.0)(0.5) + 1.0(0.2) + 3.0(0.1) + 0.1 \end{aligned}$$

$$\sigma = \frac{1}{1+e^{-z}} \quad \sigma(1.6) = \underline{0.8320}$$

$$\begin{aligned} \text{output neuron} &= [0.8320(1.2)] + [0.8176(0.7)] + 0.05 \\ &= 0.4761 \end{aligned}$$

$$\sigma(0.4761) = \boxed{0.6168} \quad \begin{matrix} > 0.5 \\ \sigma = < 0.5 = 1 \end{matrix}$$

$y_{\text{Pred}} = 1$

★ 3 important things

Activation function



Forward Propagation

Loss function



Optimizer



Backward Propagation

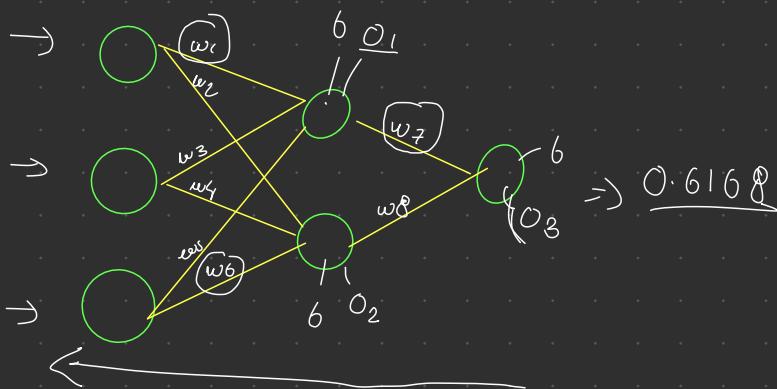
Backward propagation

$$\hat{y} = 0.6168$$

$$y = 1$$

$$\rightarrow \text{Log loss} = [y \times \log(\hat{y}) + (1-y) \times \log(1-\hat{y})]$$

$$\Rightarrow 0.494$$



$$\text{loss} = 0.494 \downarrow \downarrow$$

$$w_{\text{new}} = w_{\text{old}} - n \left[\frac{\partial L}{\partial w_{\text{old}}} \right] \rightarrow \left(\frac{\partial L}{\partial w_{\text{old}}} = \frac{\partial L}{\partial o_3} \times \frac{\partial o_3}{\partial w_7 \text{ old}} \right)$$

$$w_7 \text{ new } = w_7 \text{ old } - n \left(\quad \right)$$

$$w_1 \text{ new } = w_1 \text{ old } - n \left(\quad \right)$$

$$w_6 \text{ new } = w_6 \text{ old } - n \left(\quad \right) \rightarrow \frac{\partial L}{\partial w_{\text{old}}} = \frac{\partial L}{\partial o_3} \times \frac{\partial o_3}{\partial o_2} \times \frac{\partial o_2}{\partial w_{\text{old}}}$$



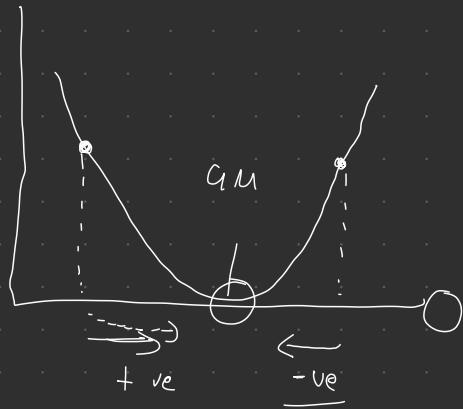
★ chain rule of Derivative ★

$$(w_1, w_2, w_3, w_4, w_5, w_6) \Rightarrow 0.432 \text{ loss}$$

↑
changed

$$\Rightarrow 0.349 - \text{loss}$$

After this we use optimizers like Gradient Descent



1 epoch \Rightarrow 71 FP
 \downarrow | BP
 \hookrightarrow 62%

2 epoch \Rightarrow 69

3 || \Rightarrow 72

4 || \Rightarrow 84

5 || \Rightarrow 86

20th epoch \Rightarrow 92%.

Activation function

(1) Sigmoid AF $\Rightarrow \phi = \frac{1}{1 + e^{-z}}$

\hookrightarrow Vanishing Gradient Problem

What is the Vanishing Gradient Problem?

The vanishing gradient problem happens during backpropagation when gradients (the small updates we calculate to adjust weights) become so tiny that the earlier layers of a deep neural network stop learning.

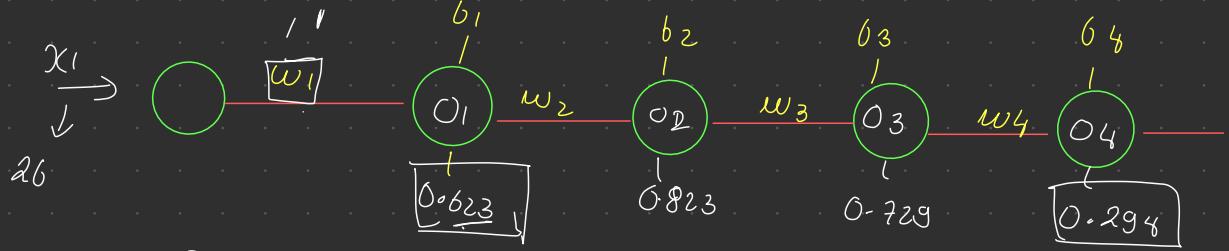
In other words:

When training deep networks, gradients are multiplied layer by layer.

If these gradients are very small (< 1), multiplying them across many layers makes them shrink toward zero.

This means the first few layers (closer to the input) never get updated properly, so the network fails to learn important low-level features.





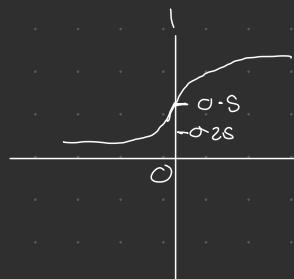
$$\sigma \left[\frac{(26x_1) + b}{1} \right] \Rightarrow \sigma = \frac{1}{1 + e^{-(26x_1) + b}}$$

0.623

100×2

$$100 \times 0.5 \\ = 50$$

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



$$\frac{\partial \text{Loss}}{\partial w_{\text{old}}} = \frac{\partial \text{Loss}}{\partial o_4} \times \frac{\partial o_4}{\partial o_3} \times \frac{\partial o_3}{\partial o_2} \times \frac{\partial o_2}{\partial o_1} \times \frac{\partial o_1}{\partial w},$$

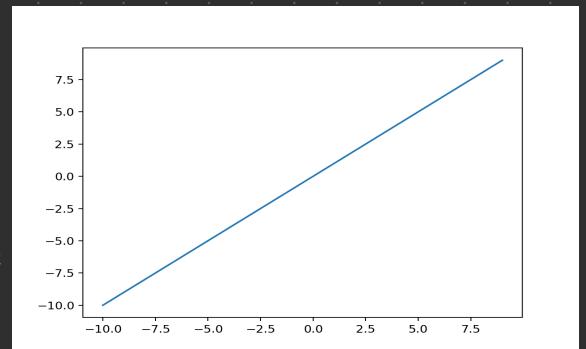
Type of Activation Functions

1. Linear Activation Function

A linear activation function is the simplest type of activation function.
It basically means:

$$f(x) = x$$

So, the output is the same as the input. The neuron doesn't transform the data it just passes it forward as it is.



Where do we use it?

In regression tasks (predicting continuous values like salary, house price, temperature, etc.).

Usually in the output layer of a neural network, because we don't want the output to be restricted between 0-1 (like sigmoid) or -1-1 (like tanh).

Example:

If you're predicting house price, you want outputs like ₹50,00,000 or ₹80,00,000.

A sigmoid function would squeeze everything between 0 and 1, which won't make sense here.

Limitations

If you use linear activation in all layers, the whole network becomes just a linear model, no matter how many layers you add.

That means it cannot capture complex, non-linear patterns in data.

That's why hidden layers use non-linear activations (like ReLU, tanh, sigmoid), but the output layer for regression can be linear.

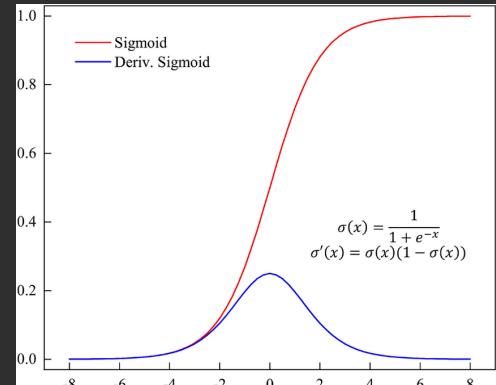
2. Sigmoid Activation Function

The sigmoid function is an S-shaped curve that squashes any real number into a range between 0 and 1.

The formula is:

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

So no matter how large or small the input, the output will always stay between 0 and 1.



Where do we use it?

In binary classification problems (e.g., predicting yes/no, disease/no disease, spam/not spam).

Usually in the output layer when you want a probability as the output.

Example:

If the sigmoid outputs 0.85, you can interpret it as 85% chance of having heart disease.

Limitations

Vanishing gradient problem: for very large or very small inputs, the gradient becomes almost 0, which slows learning.

Not used in hidden layers much nowadays (ReLU is preferred).

3. Tanh Activation Function

The tanh function (short for hyperbolic tangent) is another squashing function like sigmoid, but instead of squeezing values into 0 to 1, it squeezes them into -1 to +1.

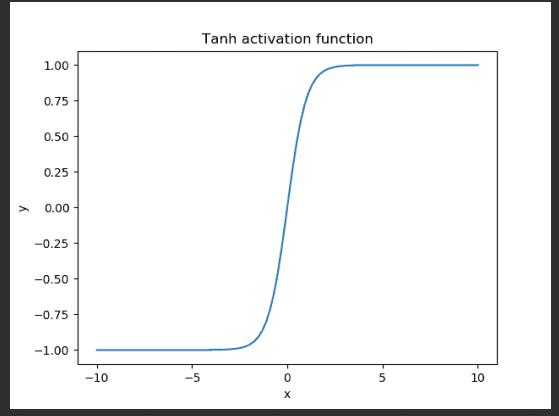
Formula :

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

Where do we use it?

- Often used in hidden layers of neural network

Useful when data has both positive and negative values because it centres the output around 0 (unlike sigmoid which is centred at 0.5)



Advantages

Outputs are zero-centered (good for optimization).

Stronger gradients than sigmoid in the range (-1,1), so learning can be faster.

Limitations

Still suffers from the vanishing gradient problem when inputs are very large (positive or negative).

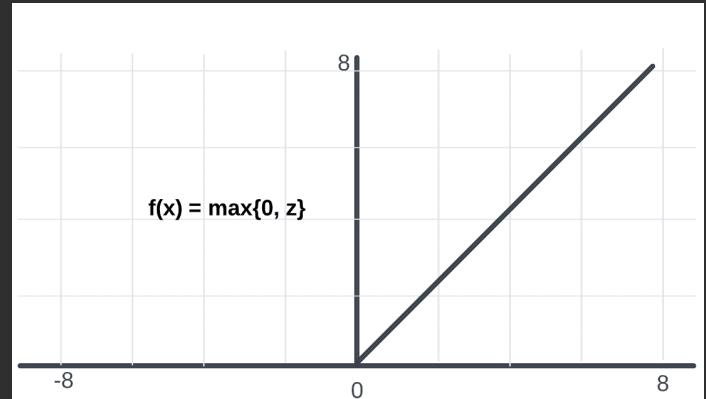
That's why in modern deep learning, ReLU is more common in hidden layers.

4. ReLU Activation Function

ReLU stands for Rectified Linear Unit.

It's super simple:

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



That means:

IF input $x < 0 \rightarrow$ output = 0

if input $x > 0 \rightarrow$ output = x

So it either passes positive values as they are or blocks negative values by turning them into 0.

Where do we use it?

- Hidden layers of almost all modern deep neural networks.
- Works really well in CNNs (Convolutional Neural Networks), image recognition, NLP, and many more tasks.

Advantages

- Very fast and simple to compute.
- Helps avoid vanishing gradient problem (better than sigmoid/tanh).
- Makes training deep networks much faster.

Limitations

- Dying ReLU problem: sometimes neurons get stuck at 0 forever if weights update badly.
- Not smooth at 0 (not differentiable there, but still works fine in practice).

5. Leaky ReLU Activation Function

It's just like ReLU, but with a small twist.

In ReLU, whenever the input is negative, the output is 0.

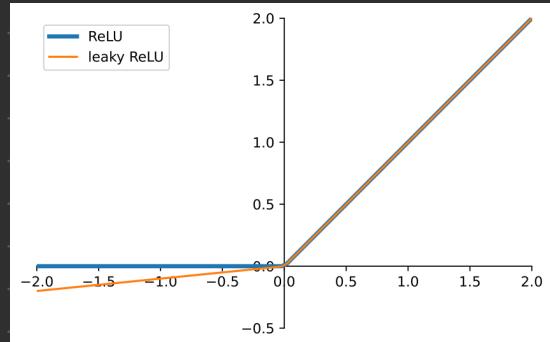
Leaky ReLU's fix:

Instead of giving 0 for negative inputs, it gives a tiny negative value (like $0.01 \times \text{input}$).

This way, the neuron is never completely dead.

Formula:

$$f(x) = \begin{cases} x & \text{if } x > 0 \\ 0.01x & \text{if } x \leq 0 \end{cases}$$



Advantages of Leaky ReLU

Fixes “Dead Neuron” Problem

- In normal ReLU, if inputs go negative, the output is always negative neurons say 0, and sometimes the neuron stops learning permanently (dead neuron).
- Leaky ReLU solves this by allowing a small negative slope, so neurons still update weights.

Computationally Simple

- Just like ReLU, the function is very easy to compute (no heavy math like exponentials in Sigmoid/Tanh).

Better Gradient Flow

- Since even negative inputs have a small gradient (e.g., 0.01), the network can continue learning, reducing the vanishing gradient issue.

Works Well in Deep Networks

- Especially useful in deep neural networks where ReLU may suffer from many dead neurons.

Limitations: small negative slope may bias results, slope value needs tuning

6. PReLU Activation Function

PReLU (Parametric Rectified Linear Unit) is an improved version of Leaky ReLU. In Leaky ReLU, the slope for negative values (like 0.01) is fixed by us. But in PReLU, that slope is learned automatically by the model during training. This makes it more flexible and adaptive.

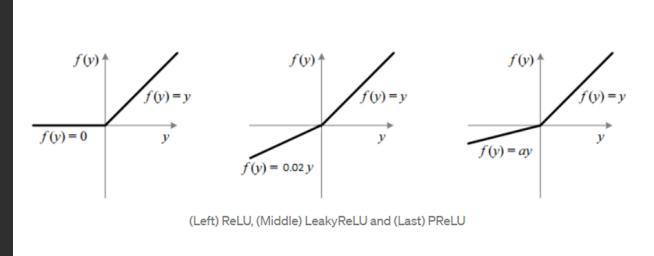
The formula is same as Leaky relu

Intuition (Easy Way)

ReLU: Negative values are killed (output = 0).

Leaky ReLU: Negative values are given a tiny leak (e.g., 0.01x).

PReLU: Instead of fixing that leak, the model says “I’ll learn the best leak slope myself.”



Advantages

1. Fixes dead neurons (like Leaky ReLU).
2. Adaptive – slope is learned, not fixed.
3. Better accuracy – often improves CNNs and deep networks.

Limitations

1. Extra parameters – slope a adds more trainable values.
2. Risk of overfitting if dataset is small.
3. Slightly more complex than plain ReLU.

7. Swish Activation Function

Swish is a smooth, non-linear activation function introduced by Google researchers.

It is Defined as :

$$f(x) = x \cdot \sigma(x)$$

Where $\sigma(x)$ is the sigmoid function
so basically: Swish = $X * \text{Sigmoid}(x)$

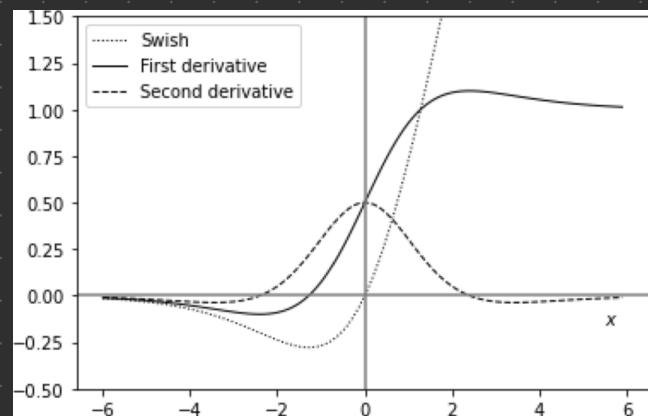
Intuition (Easy Way)

Think of it as ReLU but smoother.

For large positive inputs \rightarrow output \approx input (like ReLU).

For large negative inputs \rightarrow output is small but not strictly zero (like Leaky ReLU).

Around zero \rightarrow the curve is smooth, not sharp like ReLU.



This smoothness often makes training deep networks easier.

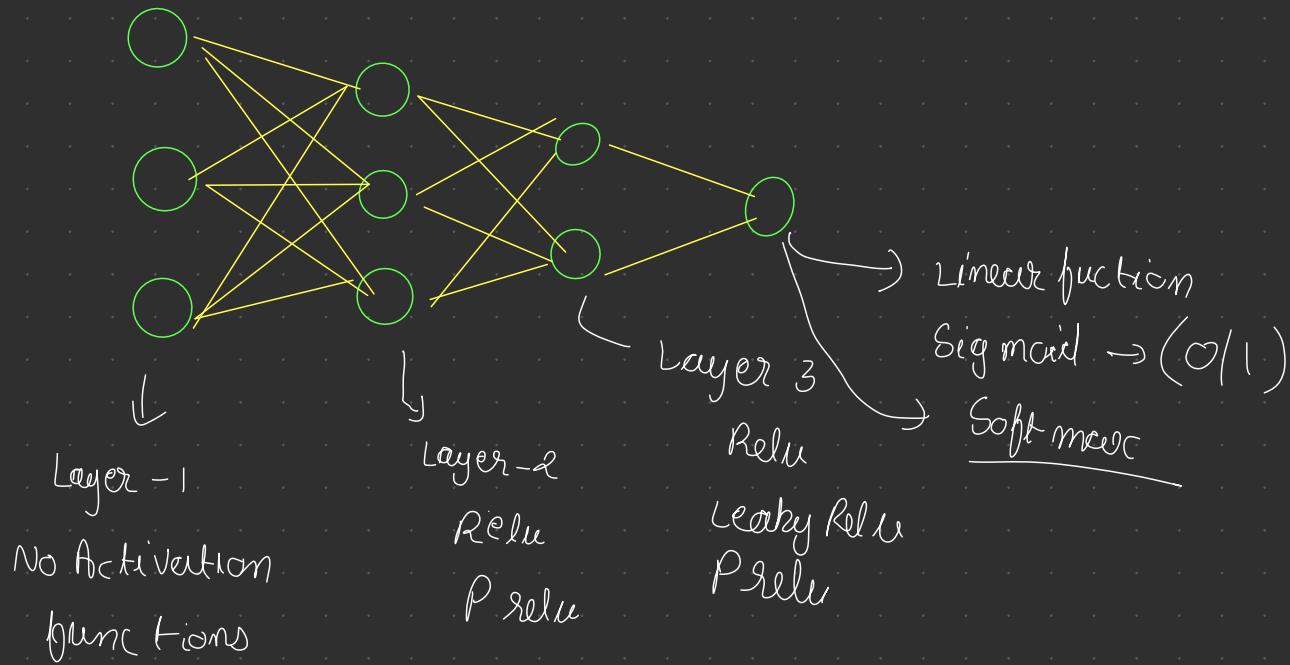
Advantages

1. Smooth curve → better gradient flow, avoids sharp jumps like ReLU.
2. Non-monotonic → can adapt better to complex patterns.
3. Works well in deep networks (often improves accuracy over ReLU).

Limitations

1. More computation (needs sigmoid).
2. Not always better than ReLU (depends on problem).
3. Slight risk of slower training compared to simple ReLU.

Usage of Activation functions



1. loss

This is the training loss.

It tells you how wrong your model's predictions are on the training set, according to the loss function you chose (in your case: binary cross-entropy).

Lower is better.

During training, the model tries to reduce this number by updating weights with gradient descent.

2. accuracy

This is the training accuracy.

It tells you what fraction of training samples the model is correctly predicting.

Example: If you have 100 samples and the model gets 80 correct → accuracy = 0.80 (80%).

3. val_loss

This is the validation loss.

After generalizing well or just memorizing the set (the split you gave with validation_data=...) and calculates the loss on that unseen data.

This tells you whether your model is generalizing well or just memorizing the training data.

4. val_accuracy

This is the validation accuracy.

It's the percentage of correct predictions on the validation set.

If accuracy is high but val_accuracy is low → your model is probably overfitting (memorizing training but failing on new data).

Where to pay attention

- loss & accuracy = how well model is doing on training data.
- val_loss & val_accuracy = how well model is doing on new unseen data.

Always pay more attention to validation metrics because that shows how the model will behave in the real world.

2000 → 2010

loss functions

Regression



MSE (Mean Squared error)

MAE (Mean Absolute error)

huber loss

MSLE

[40] ←

MSE

$$= \frac{1}{n} (\hat{y} - y)^2$$

[100] × 100

$$= 10000$$

classification

Binary cross entropy

Classification cross entropy

100 rows

(x₁) (x₂) (x₃) y



$$\text{epoch} = 100$$

①



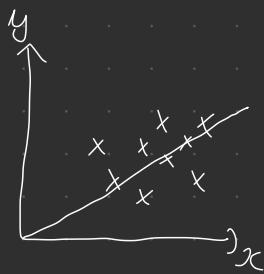
$$\text{cost function} = \frac{1}{2} \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

$$10 \times 100 \\ = \underline{1000}$$

Pros

- very common
- Penalizes the big error

$$\begin{matrix} \$K\$ & \longleftrightarrow & 100K\$ \\ & & \text{(95K\$)} \end{matrix}$$



Cons

MSE

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

Pros

- More robust to outliers
- Each error contributes equally

Cons

- Doesn't punish large errors

huber loss

Small error → MSE

Large errors → MAE

MSE

Classification loss function

one-hot encoding

Animal	cat	Dog	horse
Cat	1	0	0
Dog	0	1	0
Horse	0	0	1
Cat	1	0	0
Horse	0	0	1
Dog	0	1	0
:			

Animal

0

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

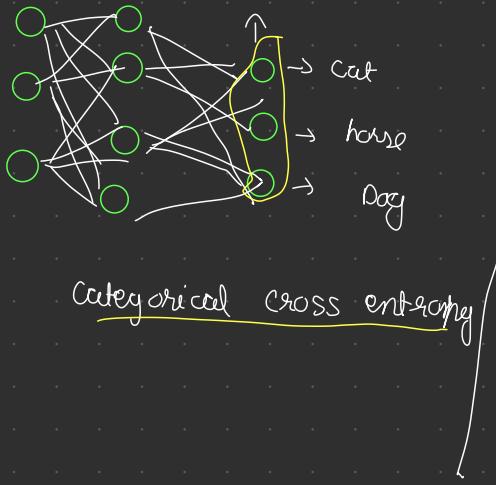
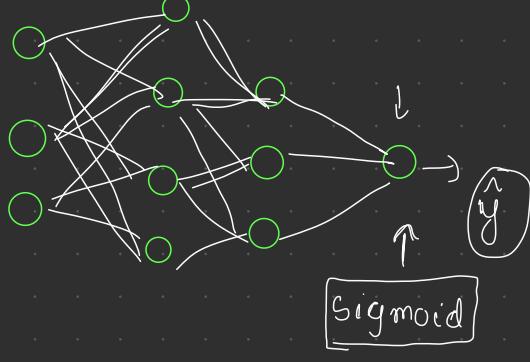
17

18

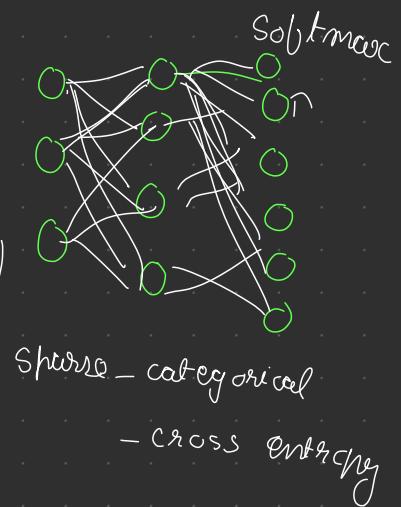
19

Softmax

Binary cross entropy



Categorical cross entropy



Binary cross entropy → Sigmoid

Loss

Activation

0, 1

Optimizers

1. Batch Gradient Descent (a.k.a. "Vanilla" Gradient Descent)

- You send all 1 lakh rows at once.
- Compute cost function on the entire dataset.
- Do 1 weight update per epoch.
- If you run 100 epochs → 100 weight updates in total.

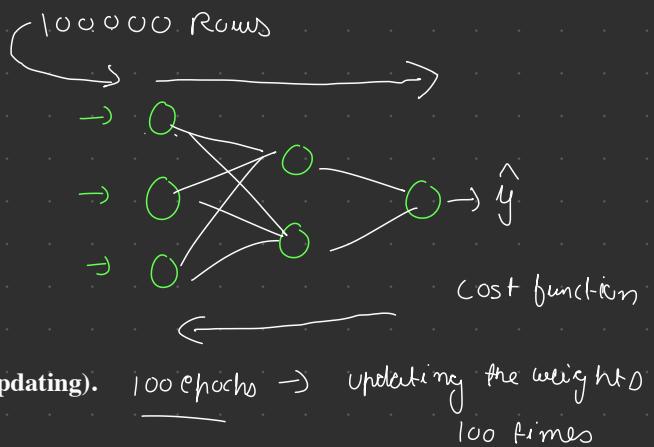
Advantages

- Smooth and stable updates (less noisy).
- Converges steadily because gradient is calculated from all data.

Disadvantages

- Slow when dataset is huge (computing gradients for all rows before updating).
- Requires a lot of RAM/VRAM to load the full dataset at once.

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



2. Stochastic Gradient Descent (SGD)

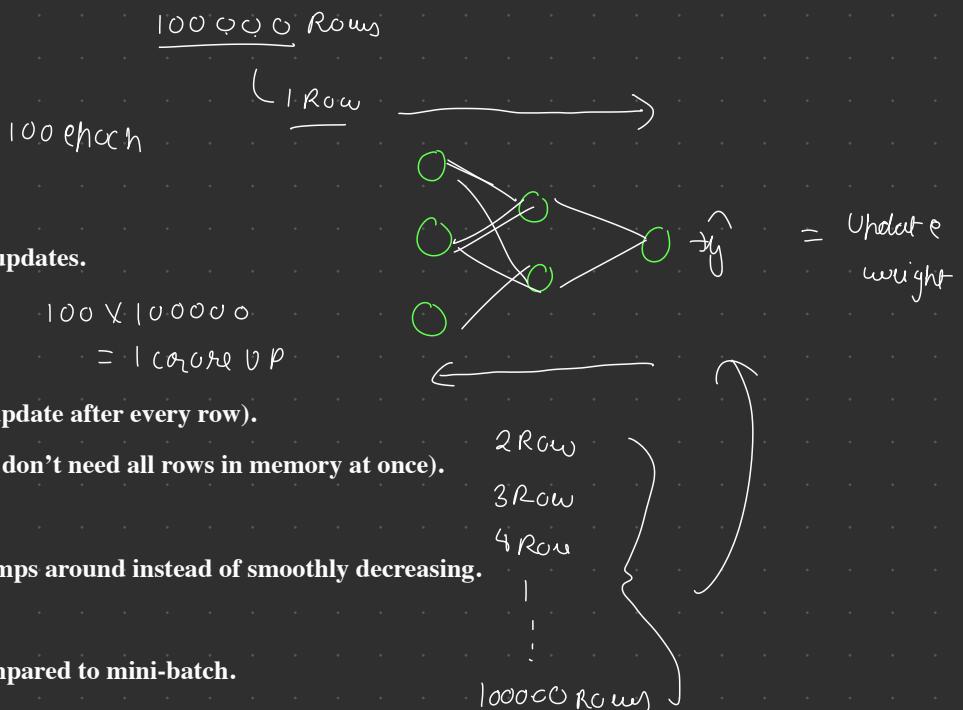
- You send 1 row at a time.
- Compute cost for that single row.
- Update weights immediately.
- With 1 lakh rows, in 1 epoch → 1 lakh updates.
- In 100 epochs → 1 crore updates.

Advantages

- Much faster to start learning (weights update after every row).
- Works well for very large datasets (you don't need all rows in memory at once).

Disadvantages

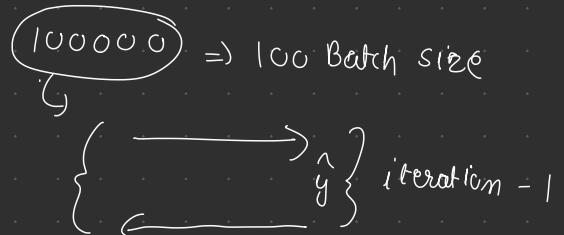
- Updates are very noisy → loss curve jumps around instead of smoothly decreasing.
- Can be unstable (harder to converge).
- Slower to reach the exact minimum compared to mini-batch.



3. The Practical Solution: Mini-Batch Gradient Descent

- A compromise between batch and SGD.
- Split data into small batches (e.g., 32, 64, 128 rows).
- Each batch → forward pass → backward pass → weight update.
- So: in 1 lakh rows, batch size 100 → 1000 updates per epoch.

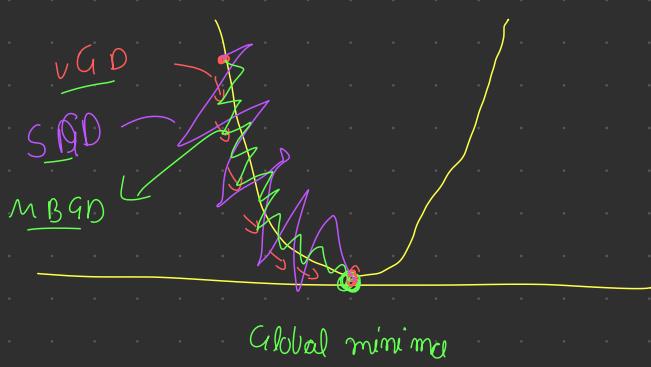
100 epochs



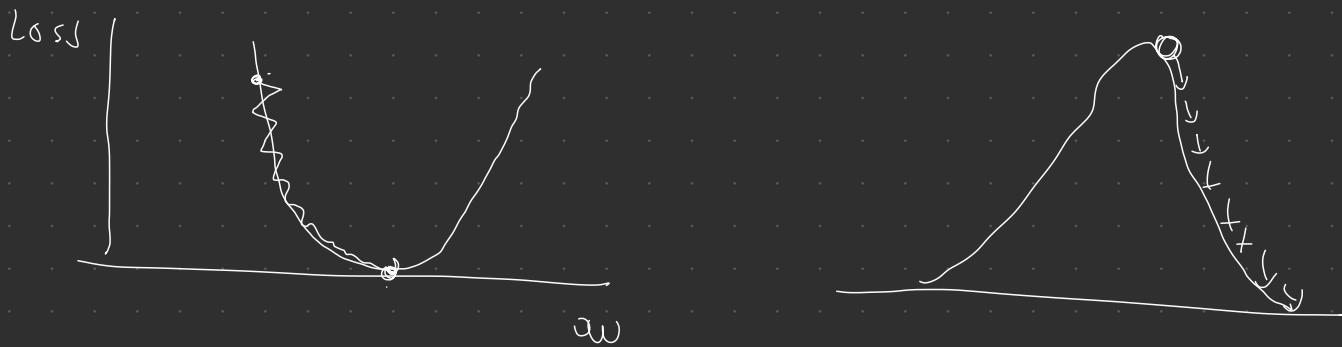
This is what almost all deep learning frameworks use today (including Keras/TensorFlow).

$$100 \times 1000 = \underline{100000}$$

$$\frac{100000}{100} = 1000$$



Momentum with SGD



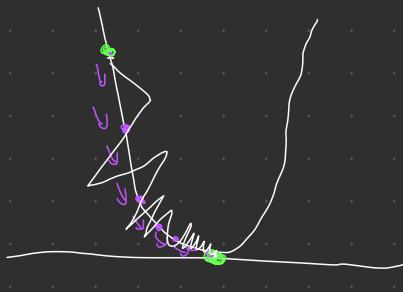
$$\text{Momentum} = \boxed{\text{mass}} \times \text{Velocity}$$

$$w_{new} = w_{old} - \eta \left(\frac{\partial L}{\partial w_{old}} \right) \Rightarrow w_{T+1} = w_T - \eta A \omega$$

$$w_{T+1} = w_T - \boxed{\sqrt{t}} \quad \boxed{0.05} \quad \Rightarrow v_t = \beta \times \boxed{0}_{t-1} + \eta A w_t$$

Adagrad (Adaptive Gradient descent)

$$w_T = w_{T-1} - \boxed{\eta} \frac{\partial L}{\partial w_{T-1}} \xrightarrow{\text{constant}} \boxed{0.01}$$

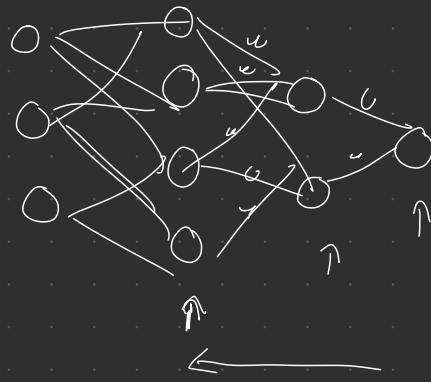


$$\omega_T = \omega_{T-1} - n \frac{\partial L}{\partial \omega_{T-1}}$$

$$m \leftarrow m' = m \rightarrow \sqrt{\alpha_T + \epsilon}$$

Epsilon

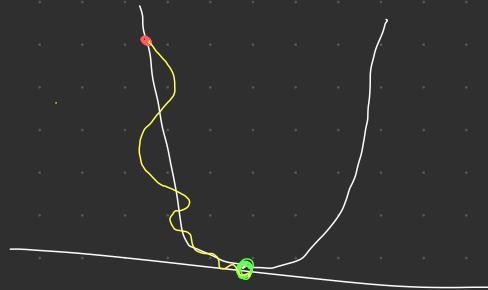
$$\alpha_T = \sum_{i=1}^t \left(\frac{\partial L}{\partial \omega_T} \right)^2 \Rightarrow$$



$$0.01 \rightarrow 0.002 \rightarrow 0.001$$

Adam optimizer

↳ Adagrad momentum



Decision Tree \rightarrow white box model

Random Forest \rightarrow Black box model

ANN \rightarrow Black box model

LR \sim white box model