



# Deep Learning Notes – Detailed, Visual & Explained

## 1. 🧠 Perceptron (1958 – Rosenblatt)



### Definition

A **Perceptron** is the **simplest neural network**, used for **binary linear classification**. It computes a weighted sum of inputs, applies an activation function, and produces a binary output.



### Working Steps

**Inputs:** (  $x_1, x_2, \dots, x_n$  )

**Weights:** (  $w_1, w_2, \dots, w_n$  )

**Bias:** (  $b$  )

**Weighted sum:**

$$z = \sum w_i x_i + b$$

**Activation:**

$$\hat{y} = f(z)$$



### Key Notes

Only works for **linearly separable** problems

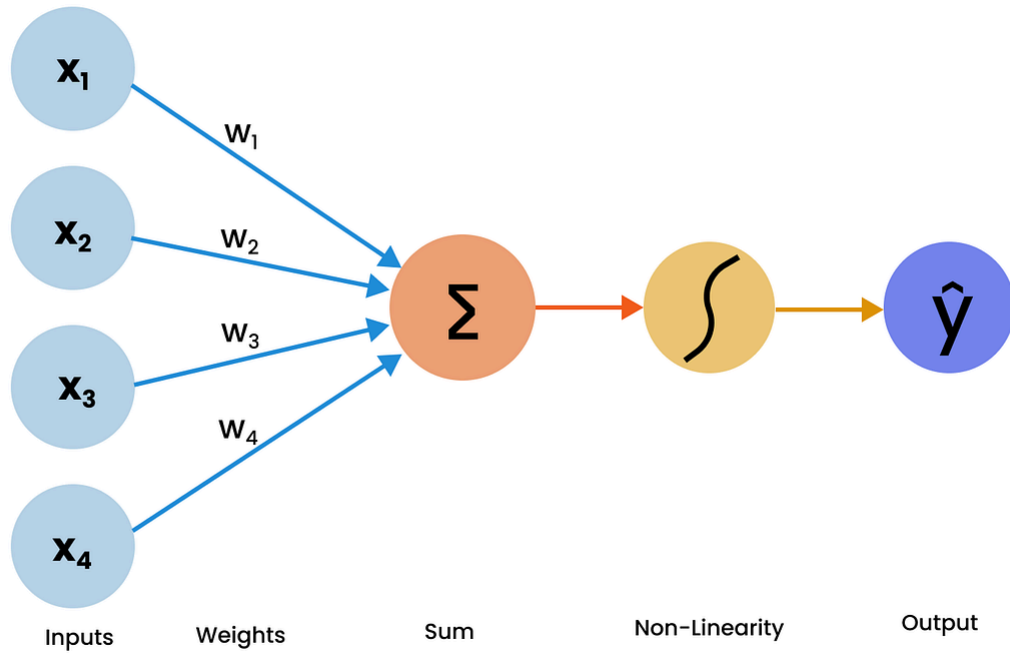
Can solve **AND/OR**, but not **XOR**



### Mental Diagram

Inputs → Weighted Sum → Bias → Activation → Output

# Perceptron



Made by Surya

## 2. 🧠 Artificial Neural Network (ANN)

### 🧠 Definition

An **ANN** is a **multi-layer neural network** that learns patterns by adjusting weights through **forward and backward propagation**. It can learn **non-linear relationships** between input and output.

### 🏗️ Structure

**Input Layer:** Receives features (pixels, text, etc.)

**Hidden Layers:** Learn intermediate patterns

**Output Layer:** Generates prediction or label

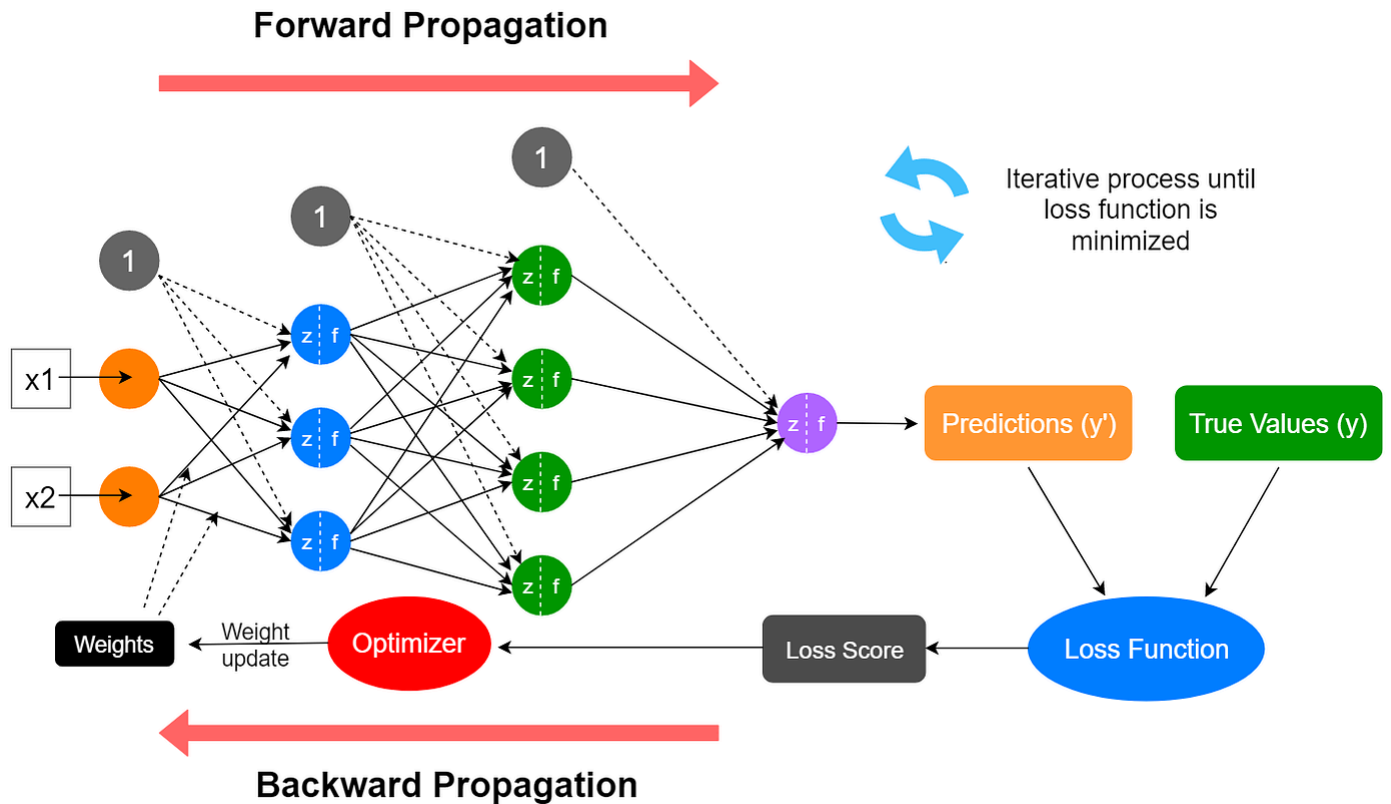
### 🔄 Learning Process

**Forward Propagation**

**Loss Calculation**

**Backward Propagation**

**Weight Update** (via Gradient Descent)



## 🚀 Why Use ANN?

Learns **non-linear** and **hierarchical** patterns

Can **approximate any function** (Universal Approximation Theorem)

Works for image, text, audio, etc.

## 📈 Visual Structure

Input → Hidden Layer(s) → Output

## ▶ Forward Propagation

Pass input through each layer

Compute activations:

$$z = w \cdot x + b, \quad a = f(z)$$

## ◀ Backward Propagation

Compute **error (loss)**

Propagate error backward using **chain rule**

Compute gradients and update weights

### 3. ❌ Vanishing Gradient Problem

#### ❌ Problem

In deep networks, gradients become very small during backprop, especially with **sigmoid/tanh**. This causes early layers to **stop learning**.

#### 📊 Example

If gradient = 0.5 per layer for 10 layers:

$$(0.5)^{10} = 0.000976$$

#### ✅ Solutions

Use **ReLU** (or variants)

Apply **Batch Normalization**

Use advanced optimizers like **Adam, RMSProp**

### 4. 🧠 Activation Functions

Activation functions introduce **non-linearity**, enabling neural networks to model complex relationships.

#### ◆ Linear

**Definition:** Outputs input as-is

**Formula:**

$$f(x) = x$$

**Range:**  $(-\infty, \infty)$

**Use:** Regression

**Cons:** ▼ No non-linearity → not useful in deep networks

#### ◆ Sigmoid

**Definition:** Squashes values between 0 and 1

### Formula:

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

**Range:** ( 0, 1 )

**Use:** Binary classification

**Cons:** ▼ Vanishing gradients, not zero-centered

### ◆ Tanh

**Definition:** Like sigmoid, but output is centered at 0

### Formula:

$$f(x) = (e^x - e^{-x}) / (e^x + e^{-x})$$

**Range:** ( -1, 1 )

**Use:** Hidden layers

**Cons:** ▼ Still vanishes at extreme values

### ◆ ReLU (Rectified Linear Unit)

**Definition:** Activates only for positive inputs

### Formula:

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

**Range:** [ 0, ∞ )

**Use:** Default for hidden layers

**Cons:** ✗ Dead neurons when (  $x < 0$  )

### ◆ Leaky ReLU

**Definition:** Small negative slope for (  $x < 0$  )

### Formula:

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

**Use:** Avoids dead ReLU neurons

### ◆ PReLU (Parametric ReLU)

**Definition:** Learns the slope for (  $x < 0$  )

Formula:

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

Use: Advanced deep networks

Swish

Definition: Smooth, non-monotonic function

Formula:

$$f(x) = x * \text{sigmoid}(x)$$

Use: Deep learning tasks in NLP & vision

Activation Function Summary

Function	Range	Use Case	Pros	Cons
Linear	$(-\infty, \infty)$	Regression	Simple	▼ No non-linearity
Sigmoid	$(0, 1)$	Binary classification	Probabilistic, smooth	▼ Vanishing gradient
Tanh	$(-1, 1)$	Hidden layers	Zero-centered	▼ Still vanishes at extremes
ReLU	$[0, \infty)$	Most common	Fast, avoids vanishing	✗ Dead neurons
Leaky ReLU	$(-\infty, \infty)$	Hidden layers	Solves dead ReLU	▼ Minor performance cost
PReLU	$(-\infty, \infty)$	Large networks	Learns best slope	▼ Adds parameters
Swish	$(-\infty, \infty)$	Deep networks	Smooth, powerful	▼ Computationally heavier

## 5. 🧠 Loss Functions

Loss functions tell us **how wrong the model is**.

### Regression Losses

#### MSE (Mean Squared Error)

$$\sum (y - \hat{y})^2 / n$$

Penalizes large errors

#### MAE (Mean Absolute Error)

$$\sum |y - \hat{y}| / n$$

Robust to outliers

#### Huber Loss

Combines MSE and MAE for stability

### Classification Losses

#### Binary Cross-Entropy:

For binary output (sigmoid)

#### Categorical Cross-Entropy:

For multi-class output (softmax)

## 6. Gradient Descent Variants

Type	Description	Pros / Cons
<b>Batch GD</b>	Uses full dataset	✅ Accurate, ▼ Slow
<b>Stochastic GD</b>	Updates per sample	✅ Fast, ▼ Noisy
<b>Mini-Batch GD</b>	Small batch updates	✅ Balanced → ★ Most used

## 7. Optimizers

Optimizer	Description	Best For
<b>SGD</b>	Basic gradient descent	Small models
<b>Momentum</b>	Adds inertia to SGD	Faster convergence
<b>Adagrad</b>	Adapts learning rate per weight	Sparse features
<b>RMSProp</b>	Exponentially decaying avg of grads	RNNs
<b>Adam</b>	RMSProp + Momentum (default choice)	Deep networks

## 8. Evaluation Metrics

Metric	Meaning
<code>loss</code>	Training error
<code>val_loss</code>	Validation error (unseen data)
<code>accuracy</code>	Correct predictions / total samples
<code>val_accuracy</code>	Accuracy on validation set

## Optional Visual Notes (Mental PNG style)



## Activation Functions:

Linear : straight line  
Sigmoid : S-curve (0 to 1)  
Tanh : S-curve (-1 to 1)  
ReLU : flat 0 for  $x < 0$ , linear  $x > 0$   
Leaky ReLU : small slope for  $x < 0$ , linear  $x > 0$   
PReLU : slope learned for  $x < 0$   
Swish : smooth curve, rises slowly negative  $\rightarrow$  positive