

**POST GRADUATE
PROGRAM IN
GENERATIVE AI
AND ML**

**Deep Learning and Neural
Network Architectures**



Module Outline

**Neural Networks and
Deep Learning
Foundation**

**Tuning and Optimizing
Deep Neural Networks**

**Convolutional Neural
Networks - I**

**Convolutional Neural
Networks - II**

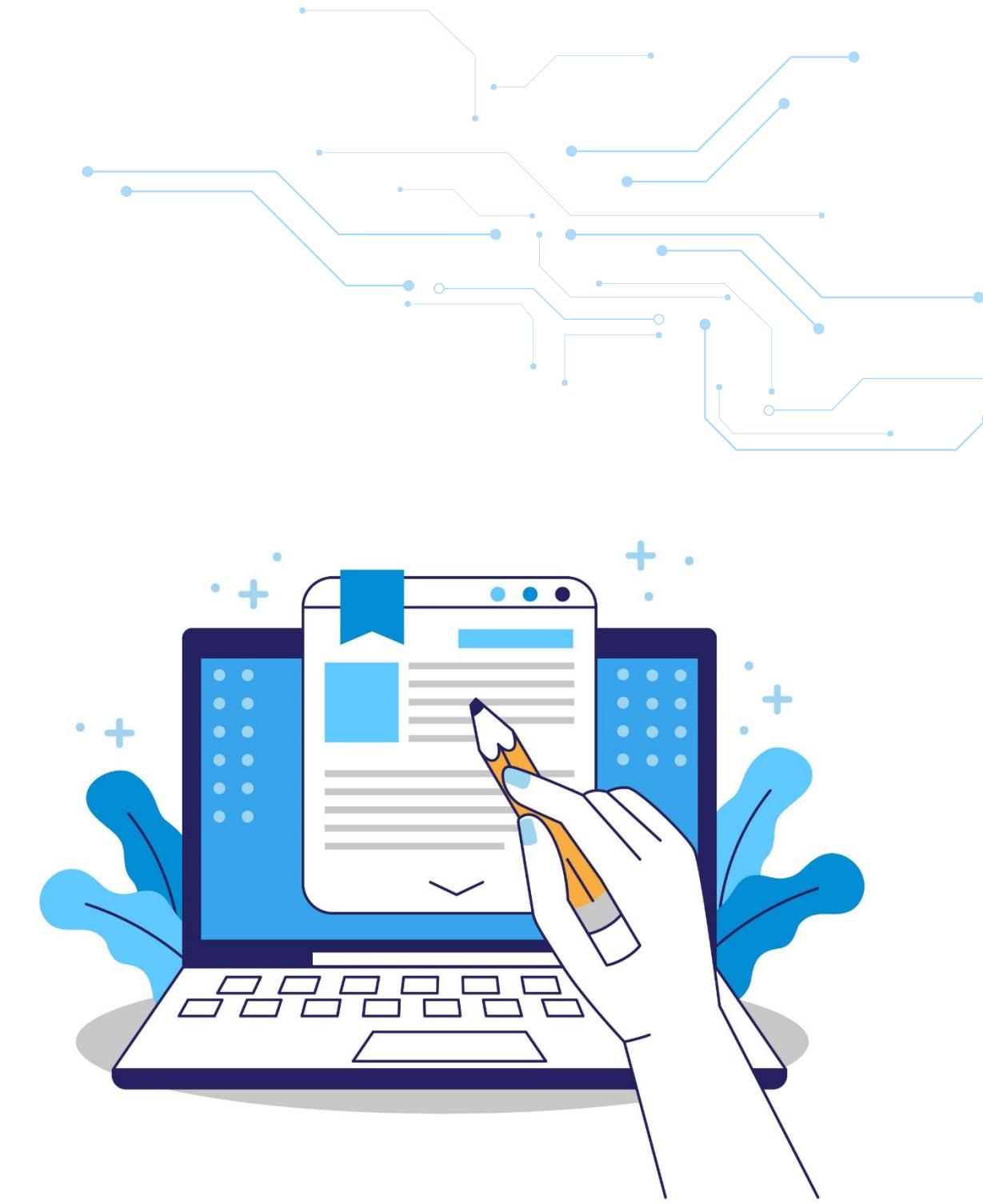
**Recurrent Neural
Networks**

**Long Short-Term
Memory (LSTM)
Networks**

Neural Networks and Deep Learning Foundation

Topics

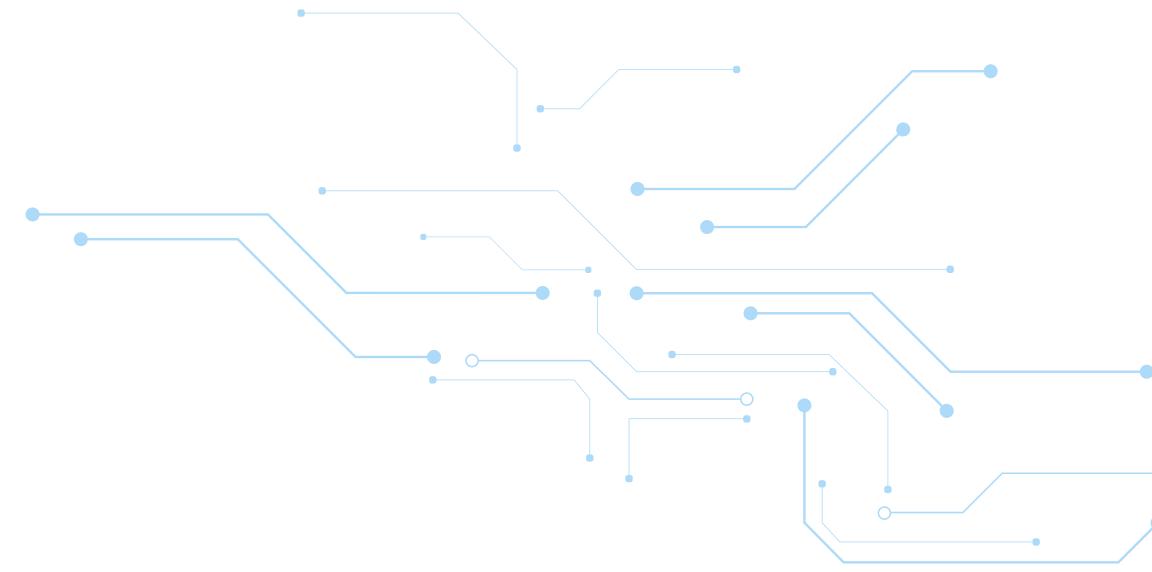
- e! Introduction to Deep Learning
- e! Biological Neuron vs. Artificial Neuron
- e! Perceptron and Multilayer Perceptron
- e! Activation, Cost and Loss Functions
- e! Forward propagation
- e! Understanding Backpropagation
Using Gradient Descent
- e! Process Of Building A Neural Network
- e! Weight initialization strategies
- e! Deep vs. Shallow Neural Networks
- e! Understanding Epochs, Batch Size,
and Learning Rate
- e! Overfitting and Underfitting
- e! Computational Graphs
- e! Monitoring Training and Validation
Loss
- e! Popular Deep Learning Frameworks
- e! Creating a simple DNN using Keras



Learning Objectives

By the end of this lesson, you will be able to:

- e! Understand deep learning basics and neural network components (neurons, MLPs, activations).
- e! Explain forward/backward propagation, cost/loss functions, and gradient descent.
- e! Build and train neural networks with proper weight initialization and tuning.
- e! Use Keras to implement and monitor basic deep learning models.



Introduction To Deep Learning

Use Case: Social Media

- e! Social media platforms like **Facebook & Pinterest** benefit heavily from Deep Learning technologies
- e! They are used for friend **recommendations**, pin **suggestions**, face **tagging** etc.

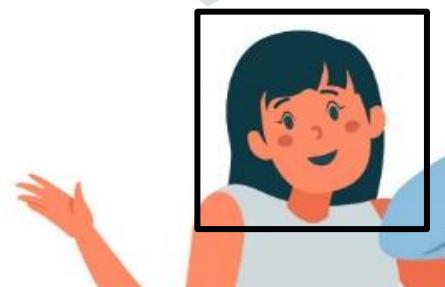


How Social Media Uses Deep Learning?

e! Example - Face tagging on Facebook:

Want to tag Susan?

- Yes - No



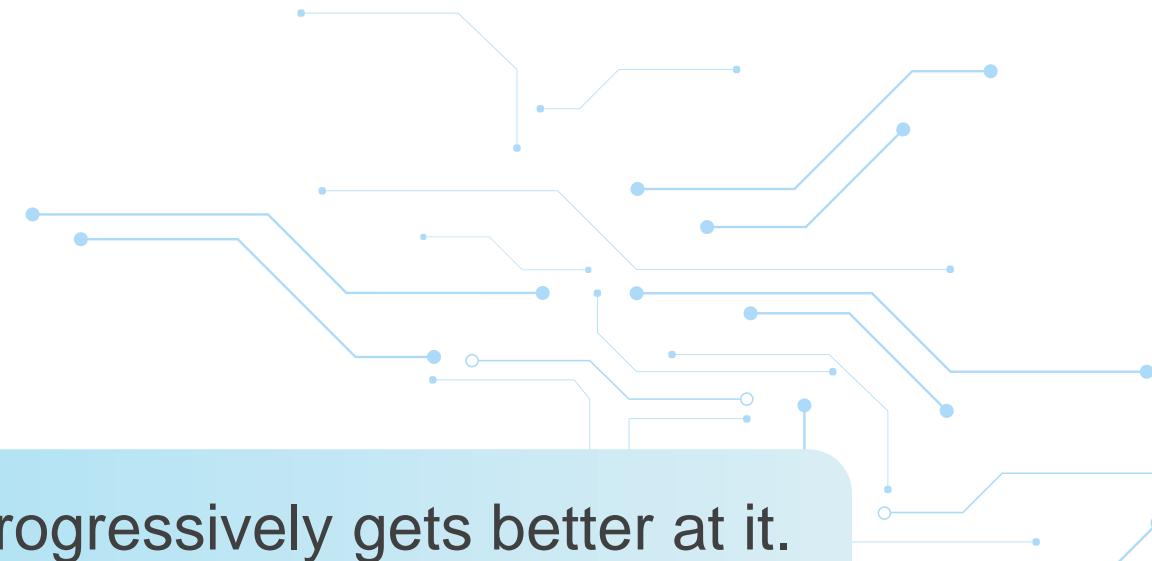
Facebook uses Deep Learning for **face recognition** in pictures uploaded by a user and then **tagging friends** from user's friend list

This information is also used to-

- e! Create **personalized feed** for the user
- e! Provide **friend suggestions**
- e! Show **relevant advertisements**

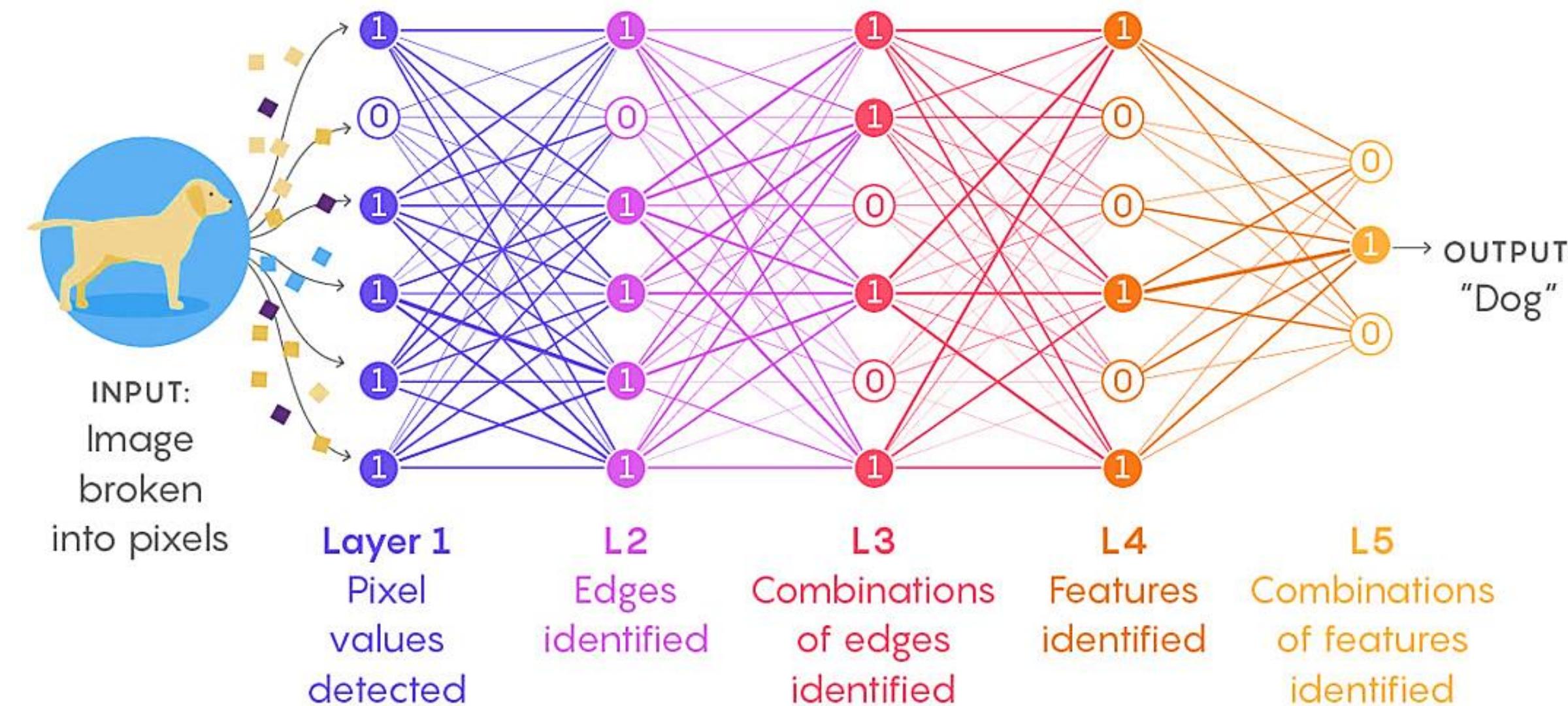


What is Deep Learning?



It is a **subset of machine learning** that takes the data, performs a function, and progressively gets better at it.

The algorithms are inspired by the structure and function of the brain, called artificial neural networks.



Why Deep Learning?

Deep learning is an Artificial Neural Network with many layers.



- e! Artificial Neural Networks (ANNs) are the core of Deep Learning
- e! ANNs are powerful, versatile and scalable
- e! ANNs find applications in-



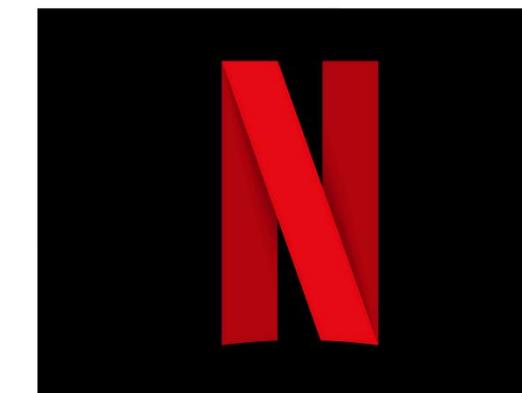
Google Images

Classifying billions of
images



Apple's Siri

Powerful Speech
Recognition



Netflix

Recommending movies
to millions of users

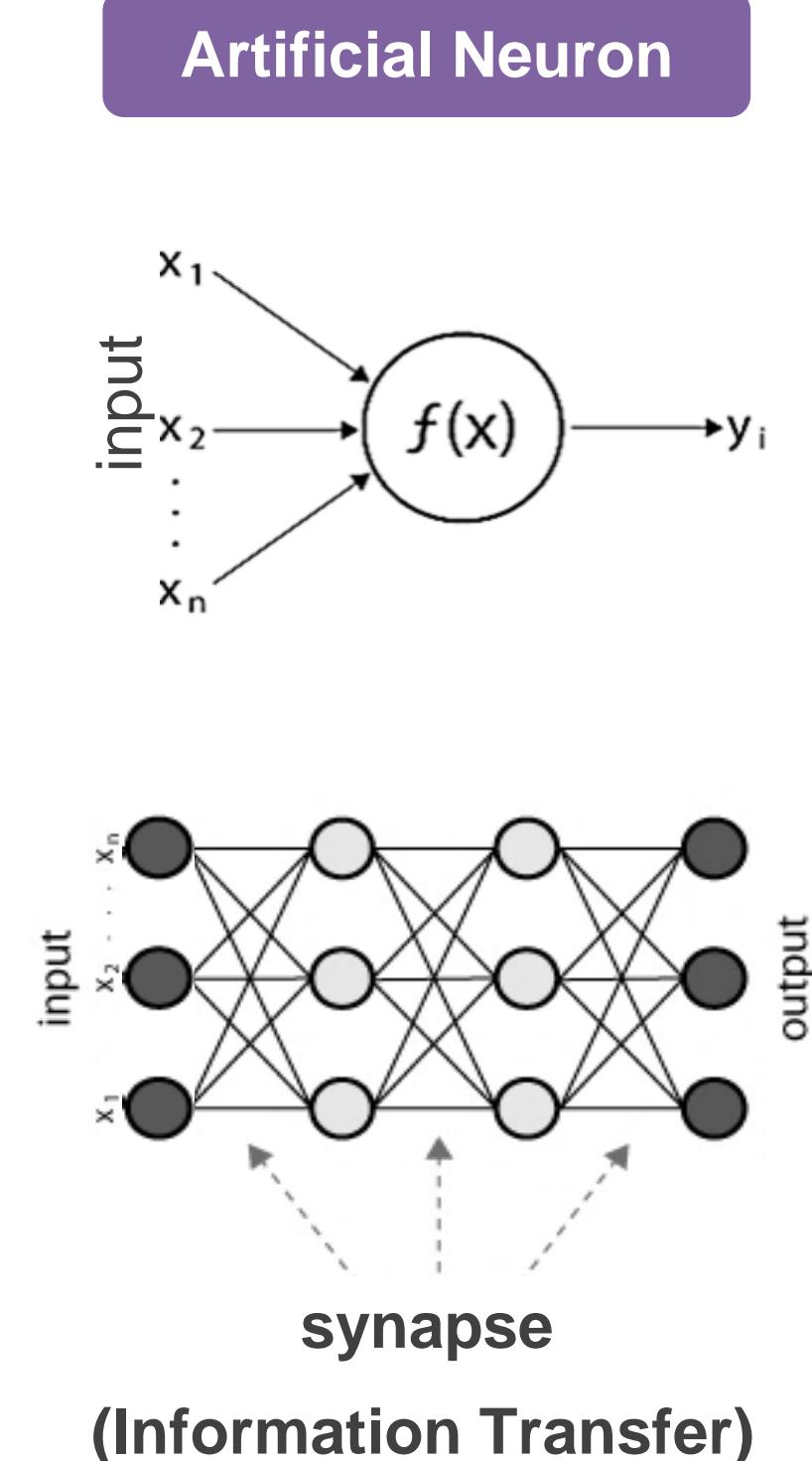
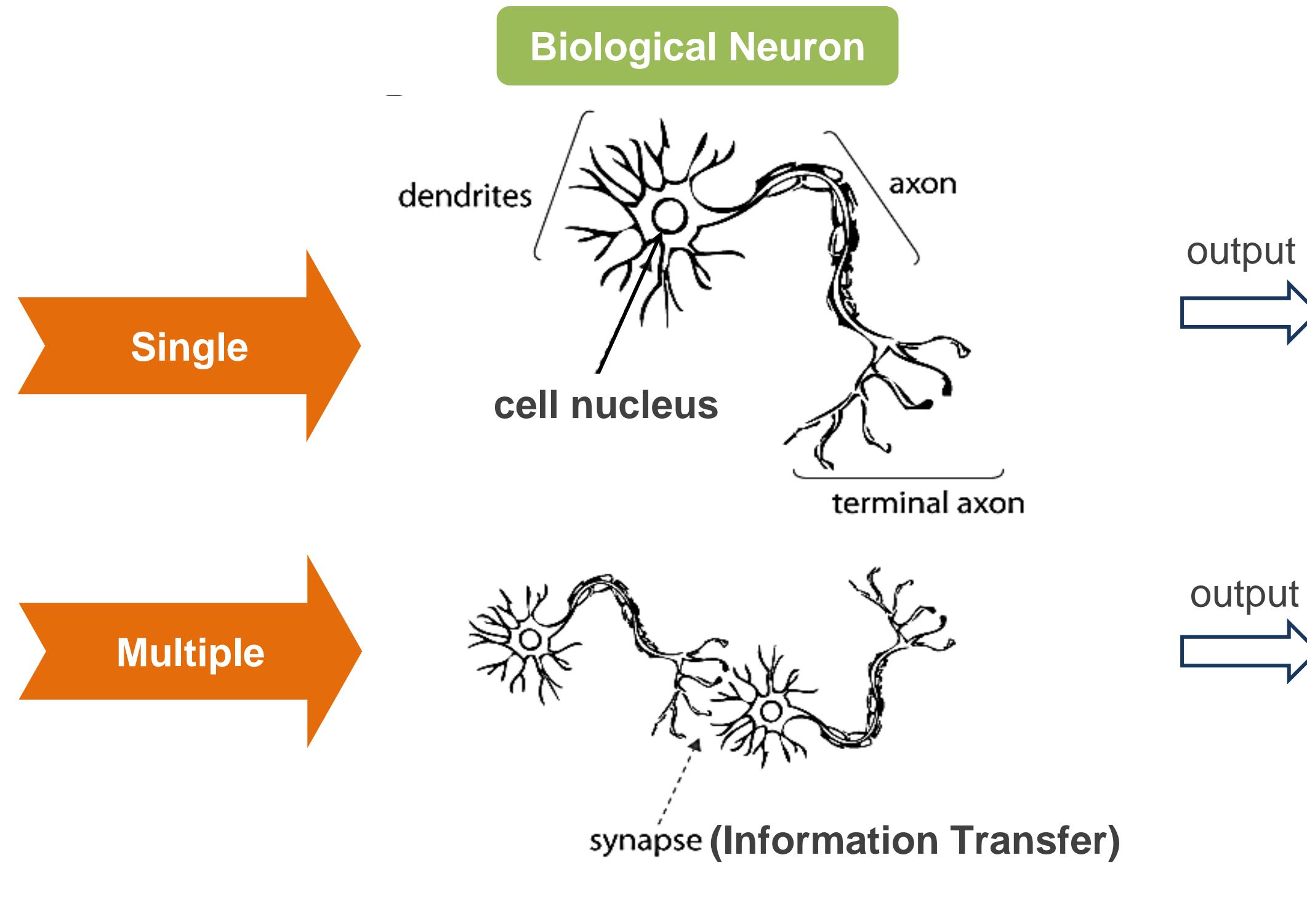


DeepMind's AlphaGo

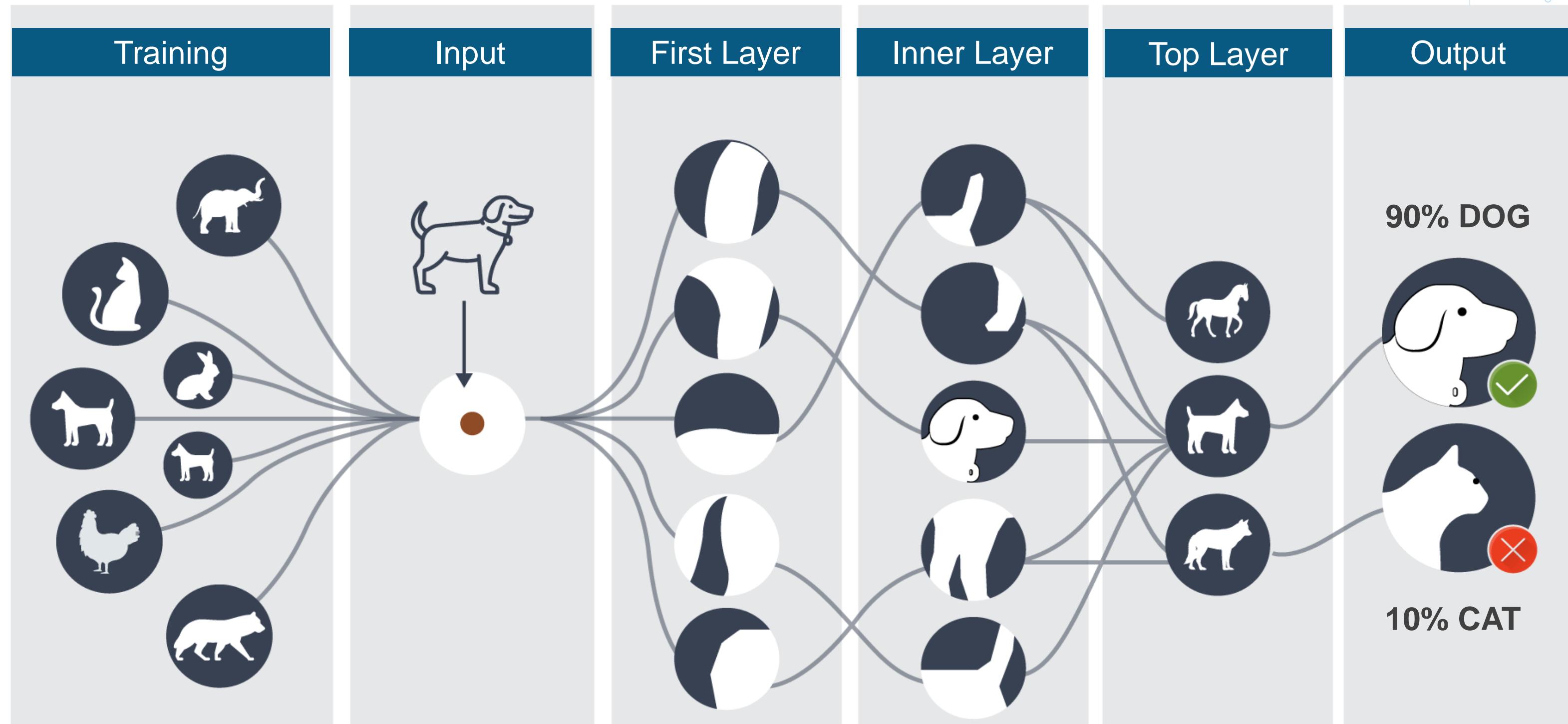
Beating the world
champion in game of Go

Biological Neuron vs. Artificial Neuron

Biological Vs. Artificial Neuron



How Does a Machine Recognize the Image?

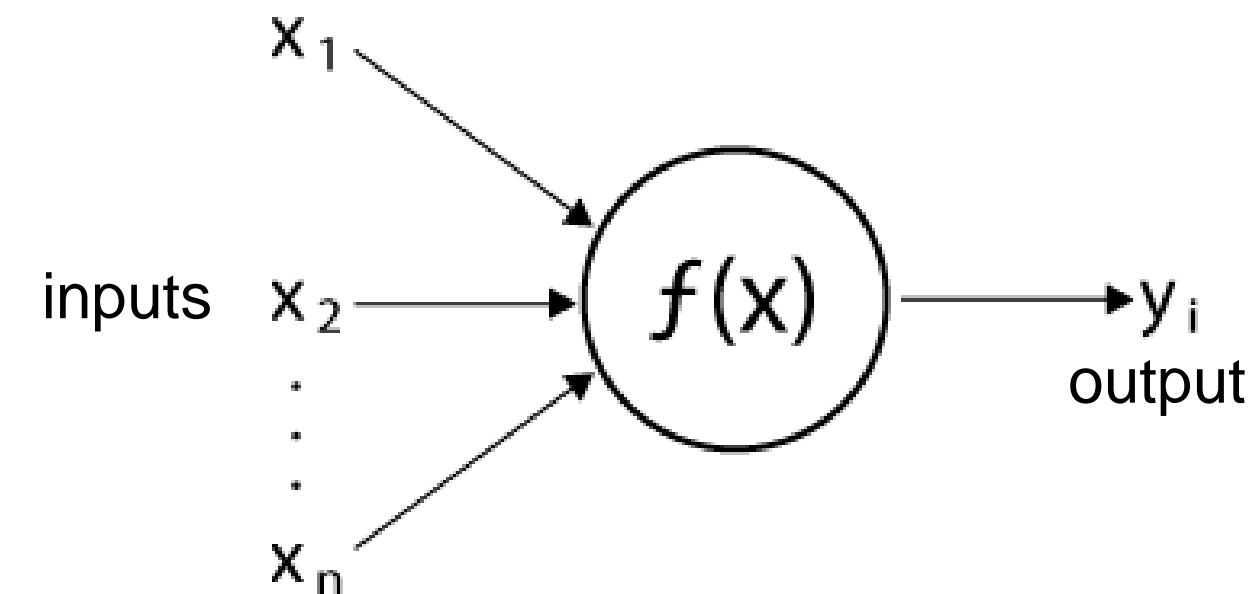


Introduction to Perceptron and Multilayer Perceptron

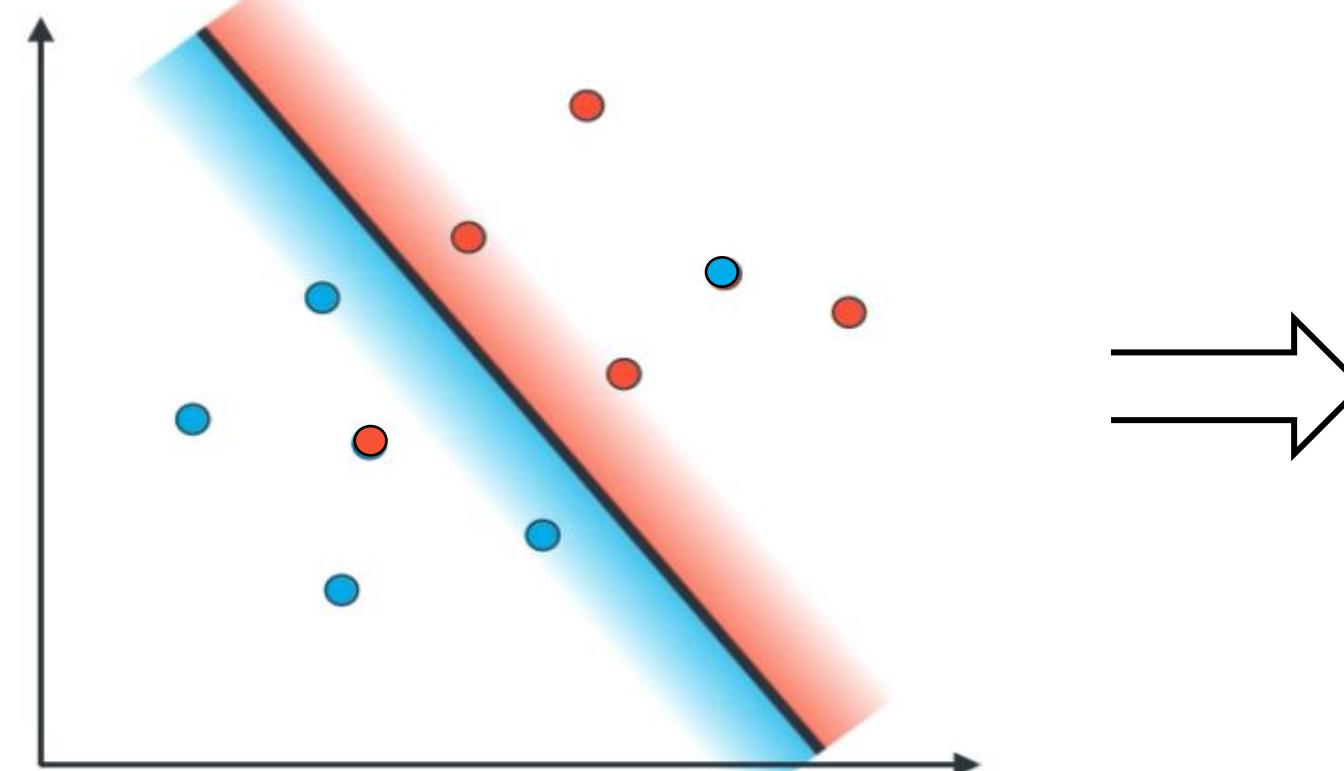
The Perceptron

A perceptron is the simplest type of artificial neural network unit, inspired by biological neurons. It takes **multiple inputs, applies weights, sums them up, and passes the result through an activation function** to produce a binary output.

Working of a Perceptron:



Limitations Of Single Layer Perceptron



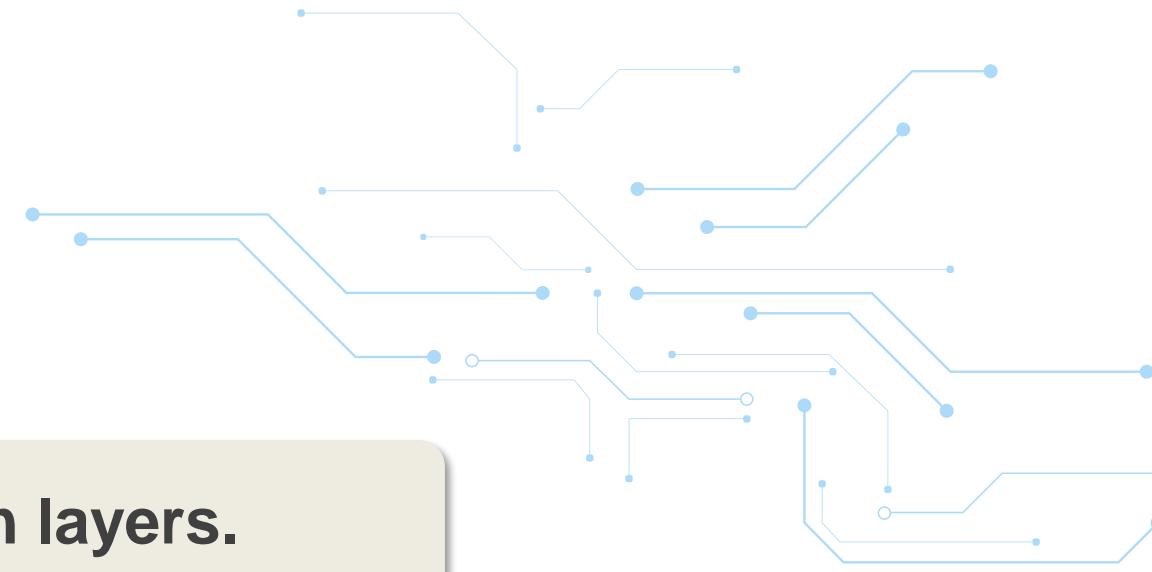
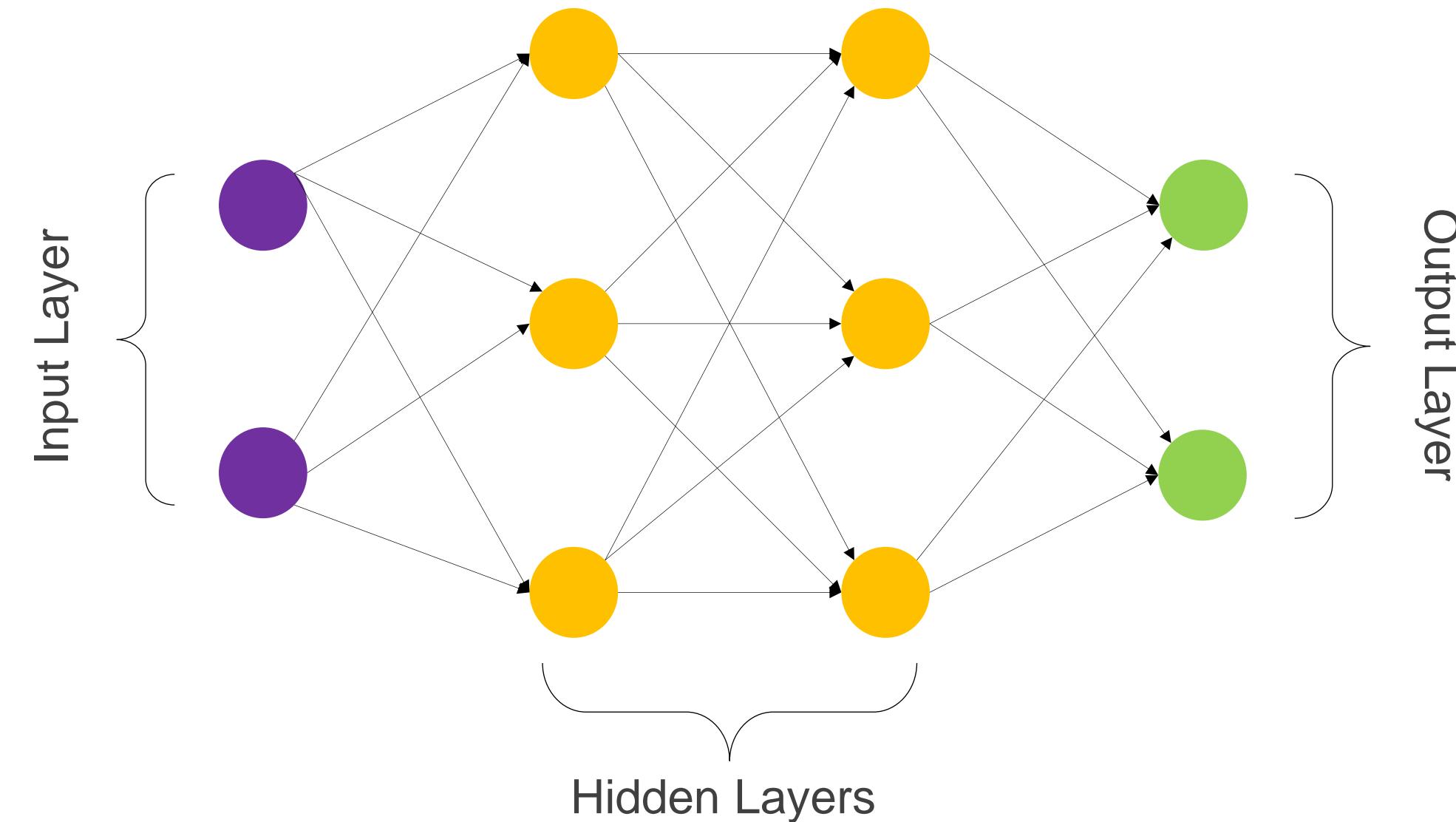
A single-layer perceptron is **not capable of separating data points with a single line**. In addition, it **can not extract complex feature dependency** and thus not good at solving complex data problems.

To solve this problem, **Multi-layer perceptron** came into picture

Multilayer Perceptron (MLP)

MLP is a artificial neural network with **more than one hidden layers**.

MLP consists of **three main layers**:



Activation, Cost and Loss Functions

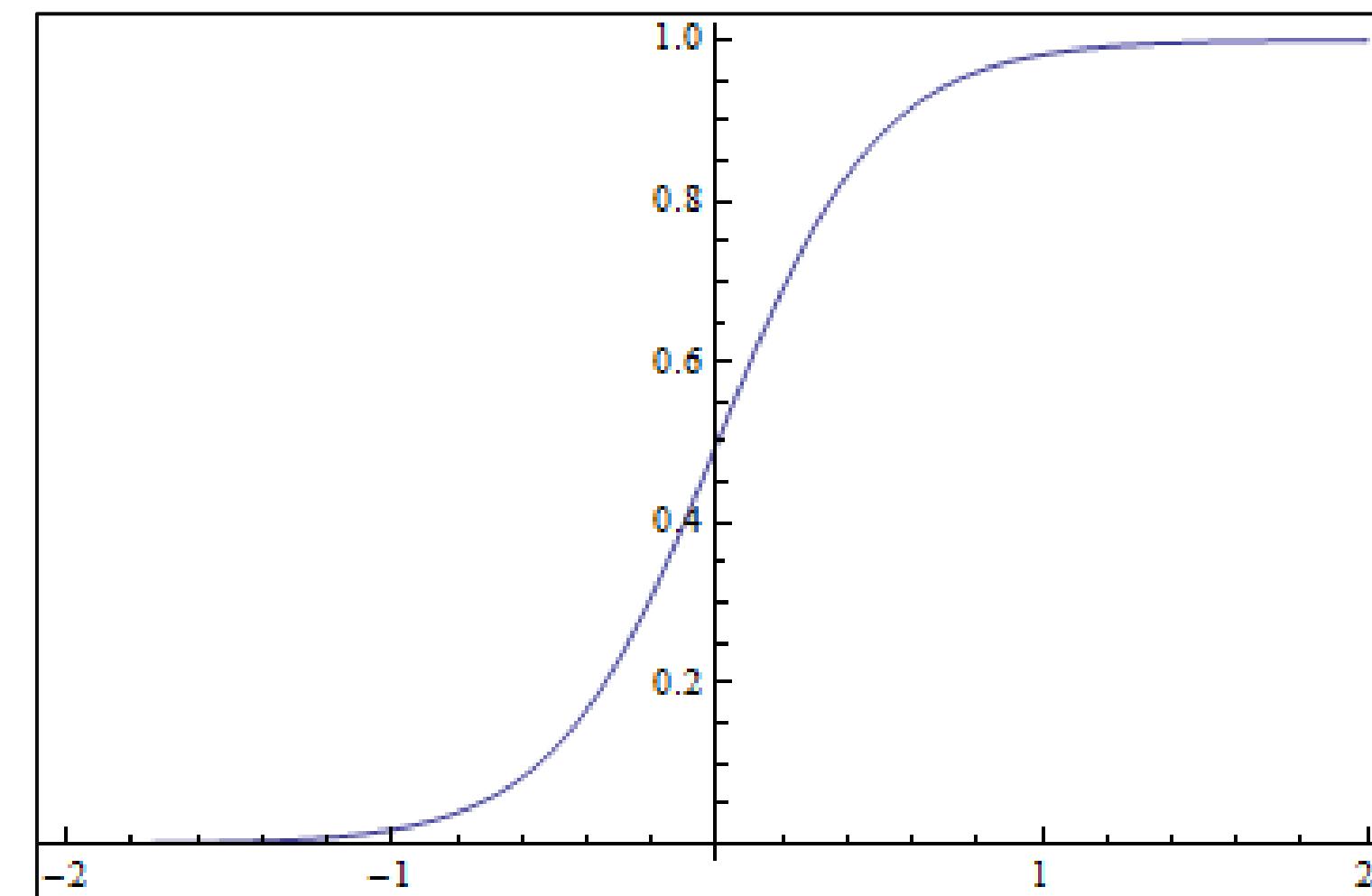
Activation Function

- e! Activation function, also known as **transfer function**, introduces non-linearity in neural networks.
- e! This non-linear transformation is introduced to **learn the complex underlying patterns** in the data.
- e! Without this function, a neuron simply resembles a linear regression.

Example: Sigmoid Function-

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

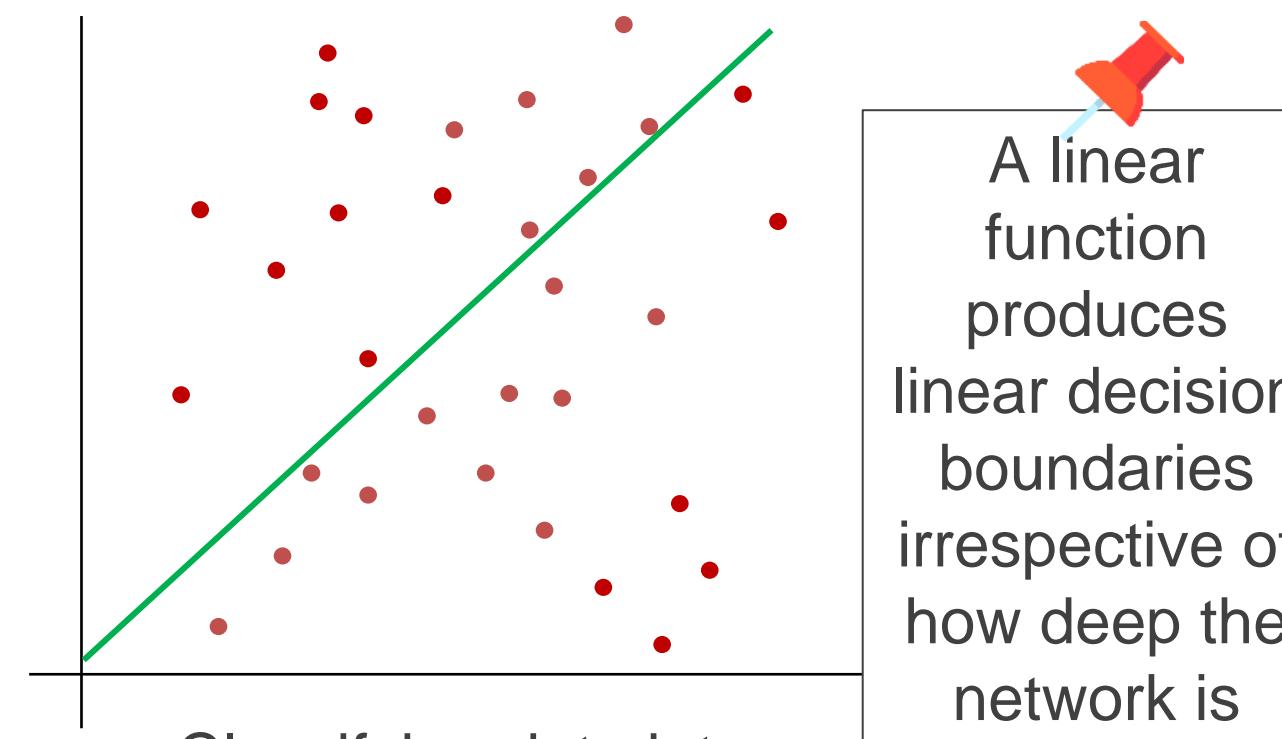
where $z = b + \sum_{i=1}^{na} x_i w_i$



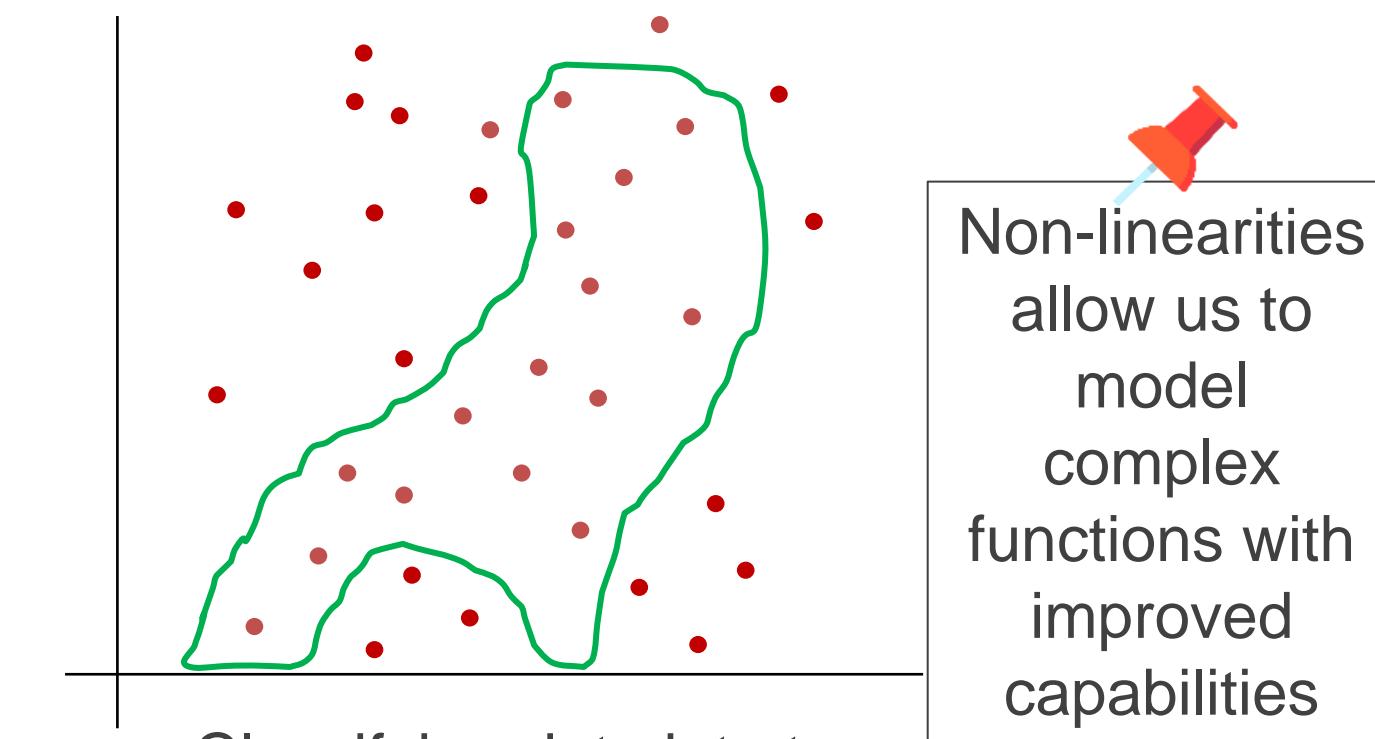
Importance Of An Activation Function

Purpose of an activation function is to **add non-linearity** into the neural network

Consider the below classification scenario:



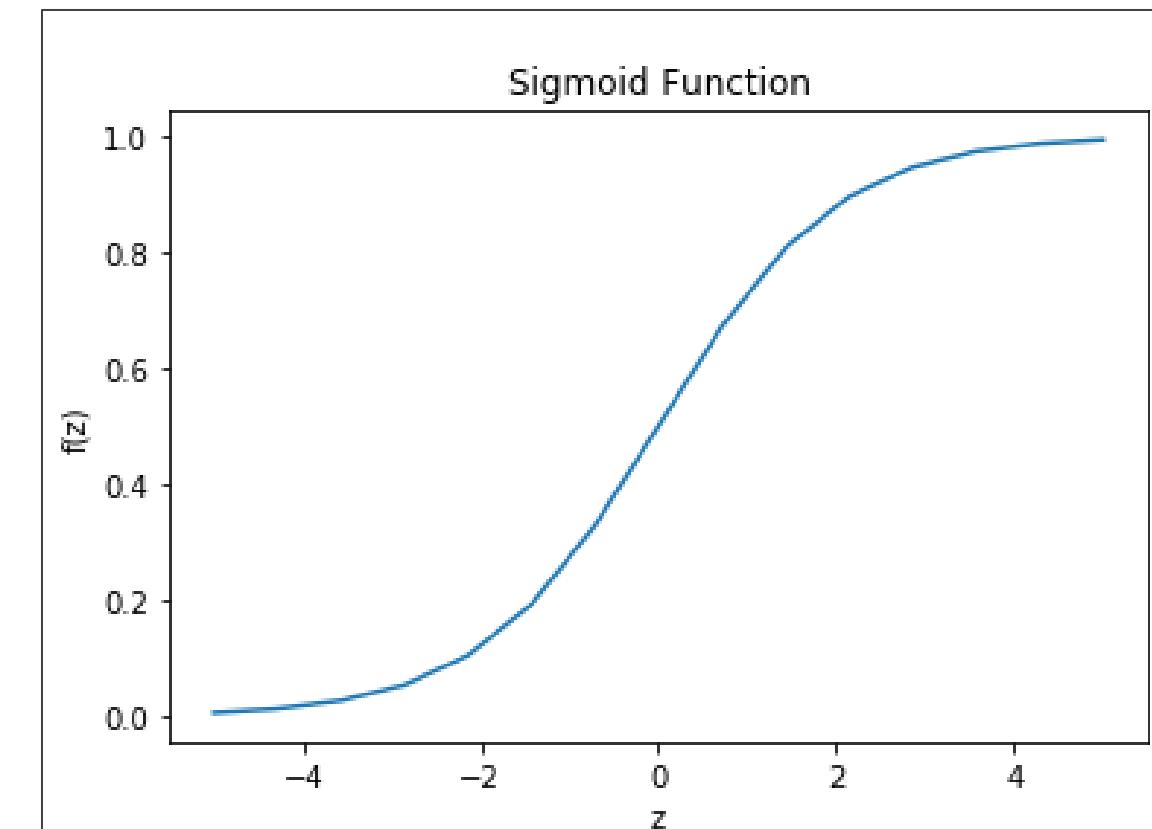
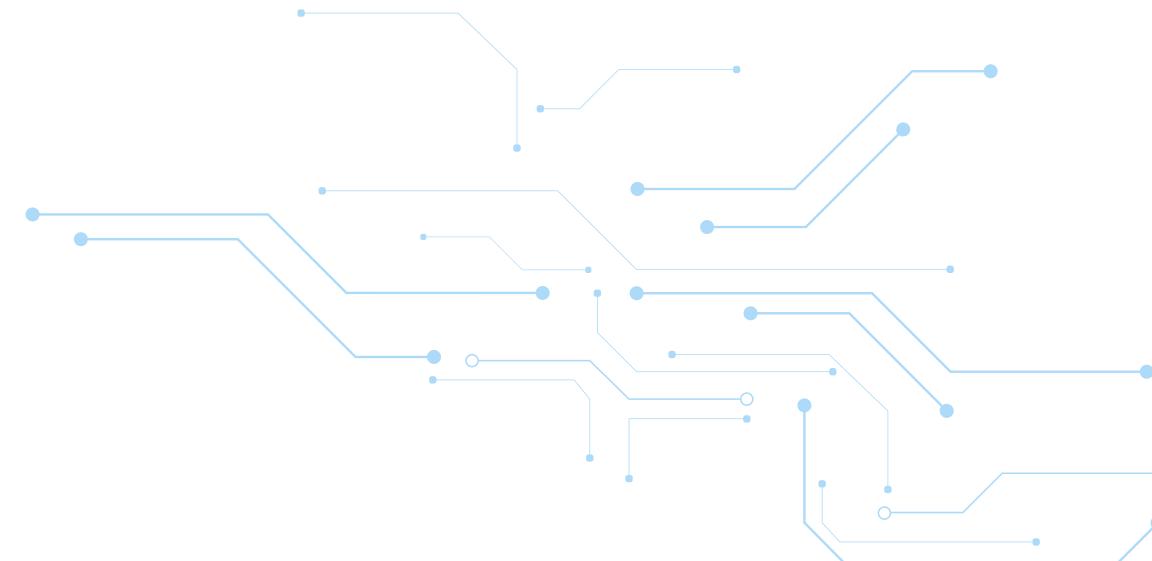
Classifying data into two labels using a linear function



Classifying data into two labels using a non-linear function

The Sigmoid Function

- e! One of the most commonly used functions.
- e! It scales all values between 0 and 1, thus used for predicting probabilities.



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

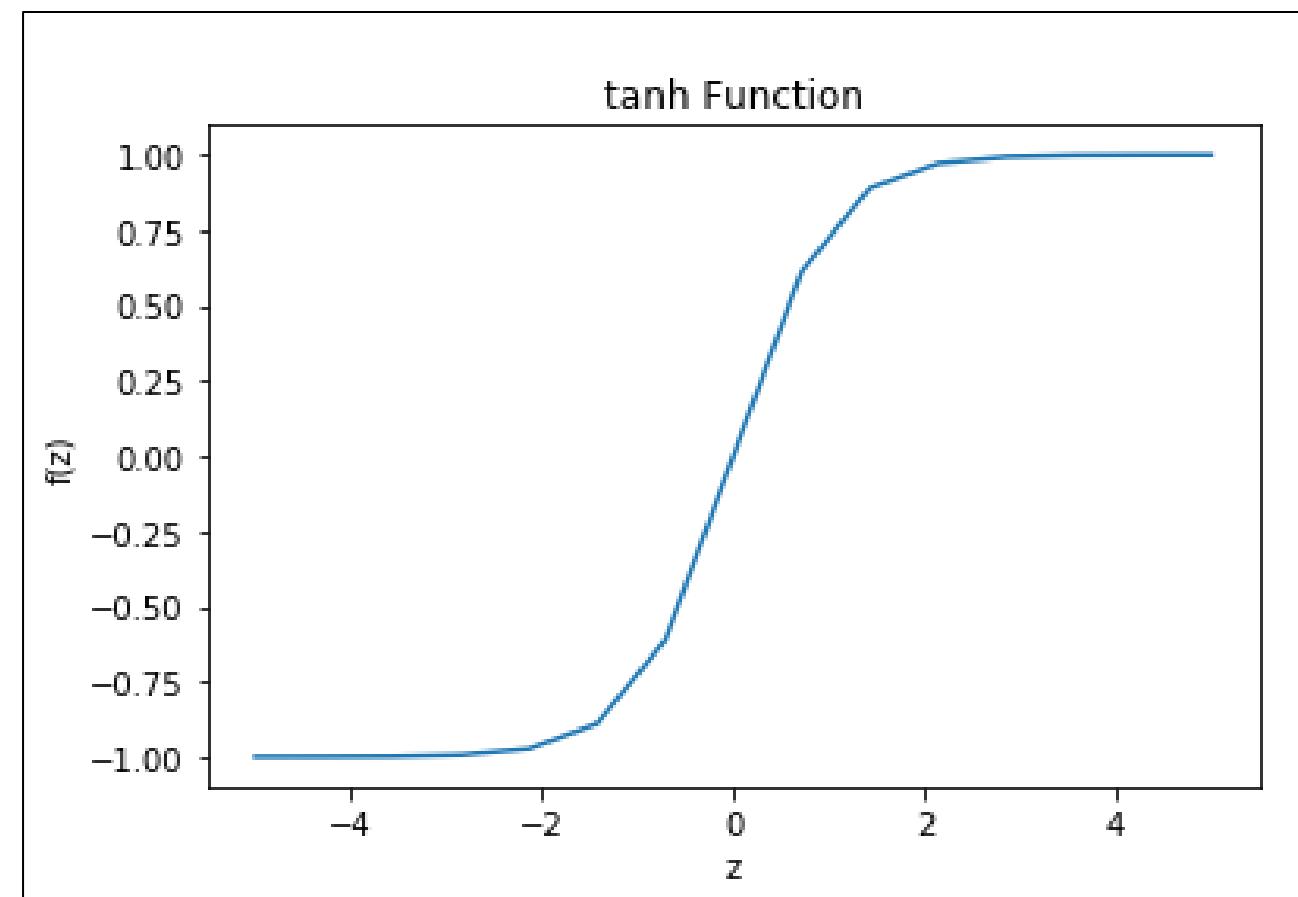


`tf.nn.sigmoid()`

It suffers from the problem of “**vanishing gradients**” which makes it unsuitable for backpropagation.

The tanh Function

- e! A hyperbolic tangent (tanh) function maps input values to outputs between (-1 , 1)
- e! This function solves the “vanishing gradients” problem to some extent

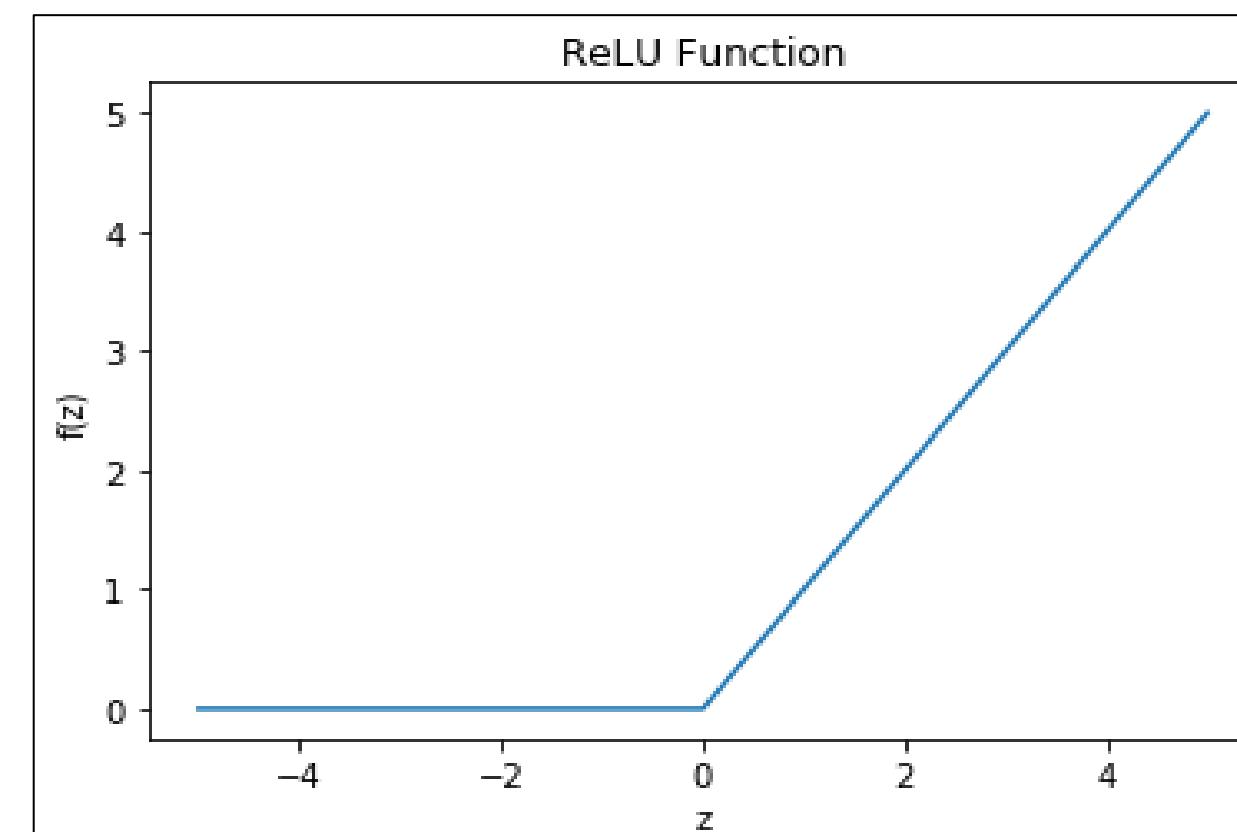


$$f(z) = \frac{(1-e^{-2z})}{(1+e^{-2z})}$$



The Rectified Linear Unit (ReLU) Function

- e! One of the most commonly used activation functions
- e! It outputs a value from 0 to ∞ as per the below equation
- e! ReLU does not saturate unlike Sigmoid & tanh functions, which saturate in both the directions. Thus, ReLU minimises the vanishing gradients problem to a great extent



$$f(z) = \begin{cases} 0, & \text{for } z < 0 \\ z, & \text{for } z \geq 0 \end{cases}$$



The Softmax Function

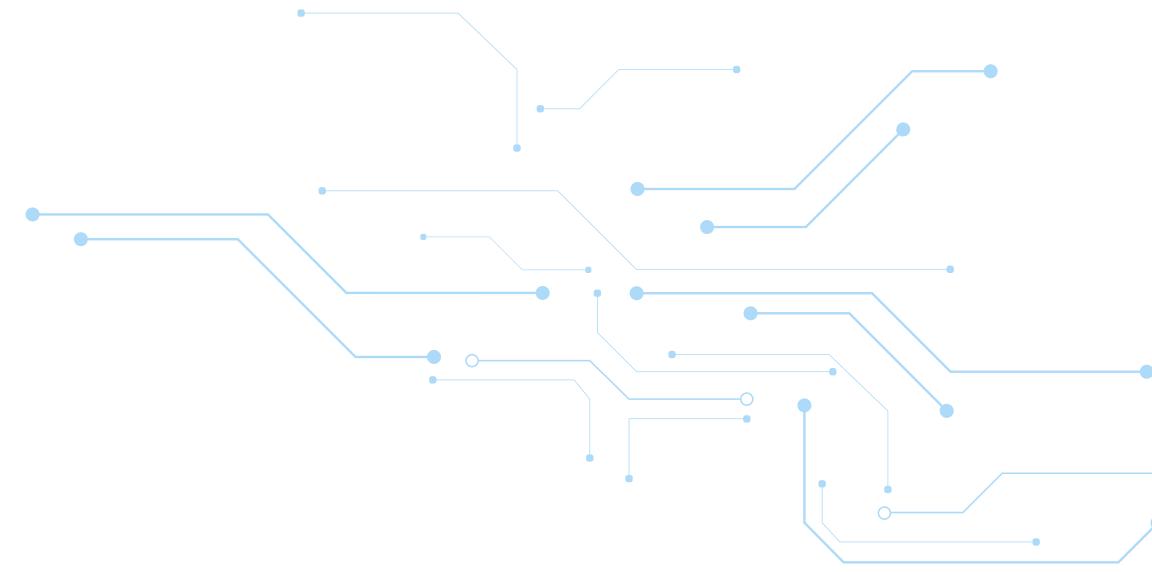
- e! It is a generalization of Sigmoid function
- e! Softmax maps outputs to values between (0,1) such that their sum is 1
- e! Thus, Softmax is applied to the final layer of a network in a multi-class classification

z_n	$f(z_n)$
-5	0
-3	0.0003
0	0.0064
2	0.0471
5	0.9462

$$f(z_n) = \frac{e^{z_n}}{\sum_{i=1}^n e^{z_n}} \text{ where, } n \text{ is the number of classes}$$



`tf.nn.softmax()`



Loss Function

e! **Mean Squared Error (MSE)**: for regression tasks

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

e! **Binary Cross-Entropy**: for binary classification

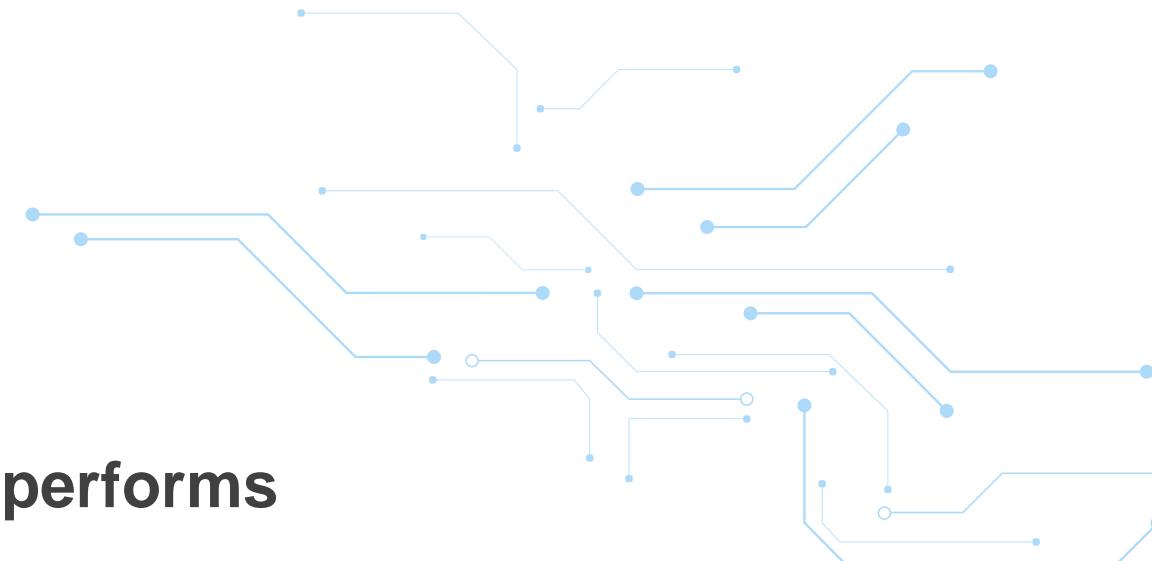
$$\text{BCE} = -\frac{1}{n} \sum_{i=1}^n y_i \cdot \log \hat{y}_i + (1 - y_i) \cdot \log (1 - \hat{y}_i)$$

e! **Categorical Cross-Entropy**: for multi-class classification.

$$\text{Loss} = - \sum_{i=1}^{\text{output size}} y_i \cdot \log \hat{y}_i$$

n: Number of training samples;
y_i: Actual Output;
(y_i): Predicted class probabilities;

Cost Function



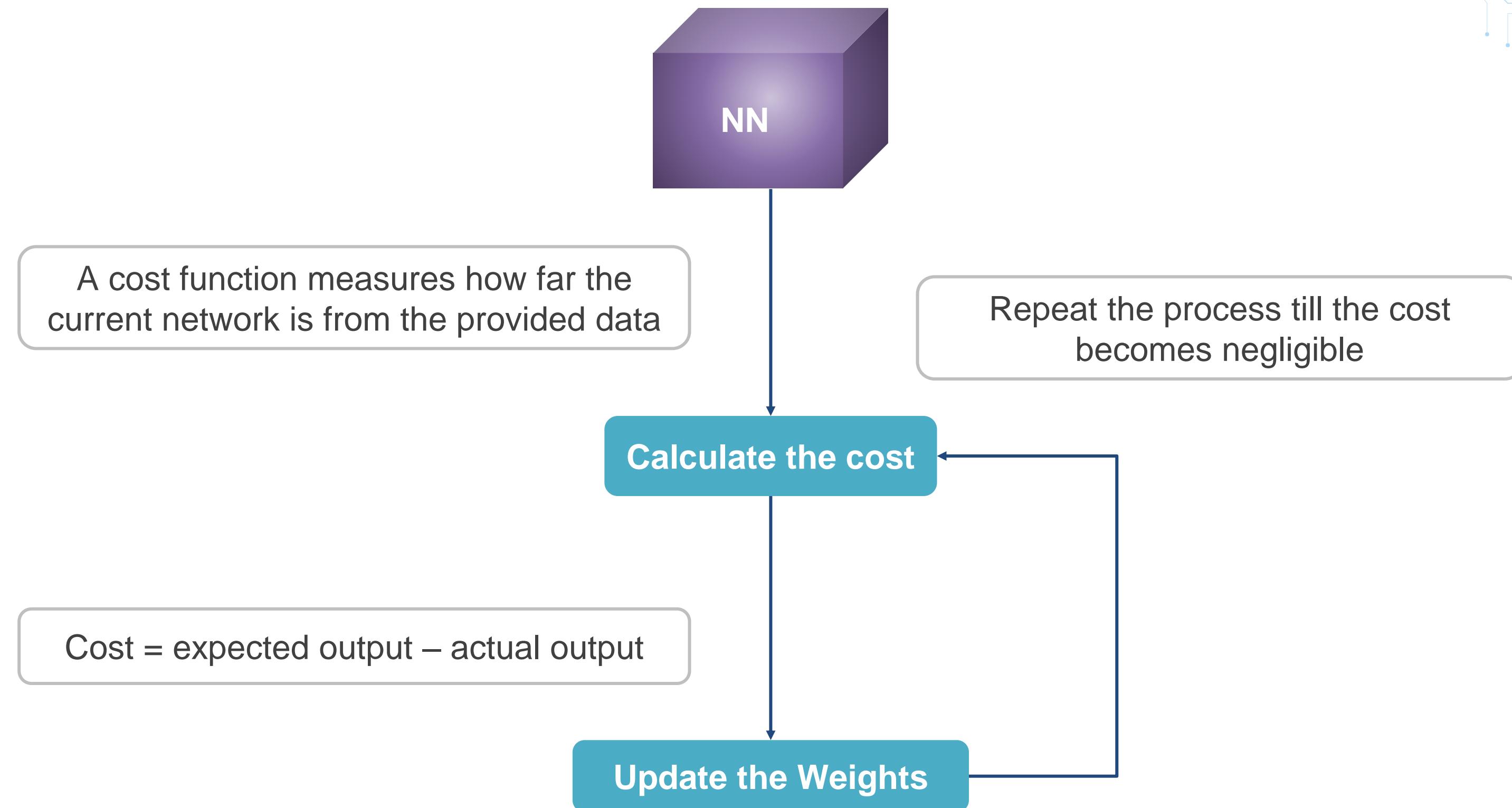
- e! A function, called Cost Function or Loss Function tells **how well a neural network performs**
- e! We will use a simple Cost Function defined as mean squared difference between the actual and the predicted output

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

- e! n : Number of training samples;
- e! y_i : Actual Output;
- e! \hat{y}_i : Predicted Output

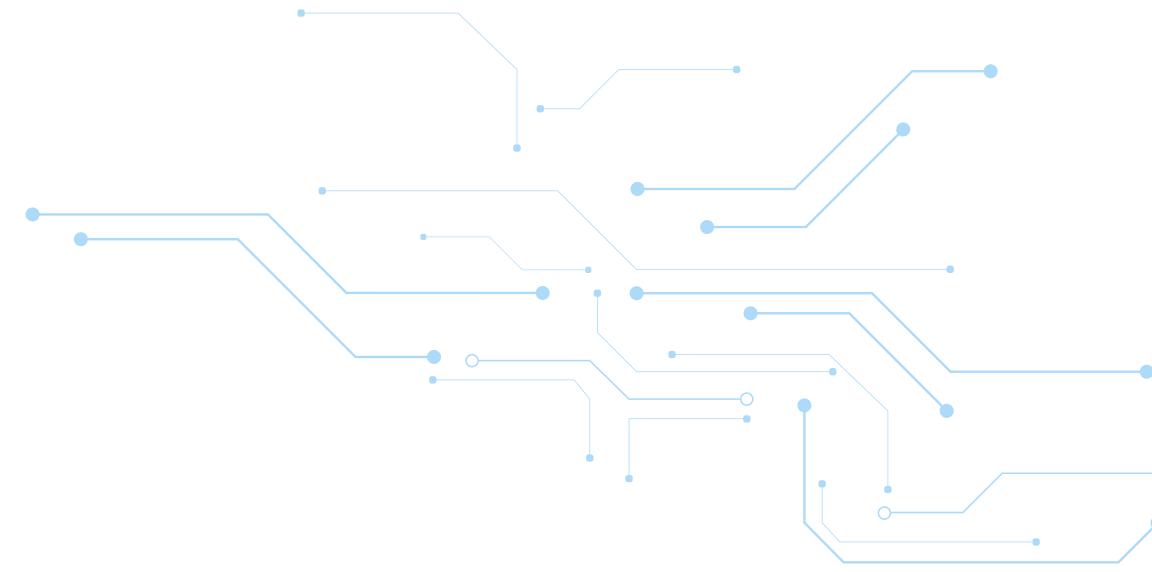
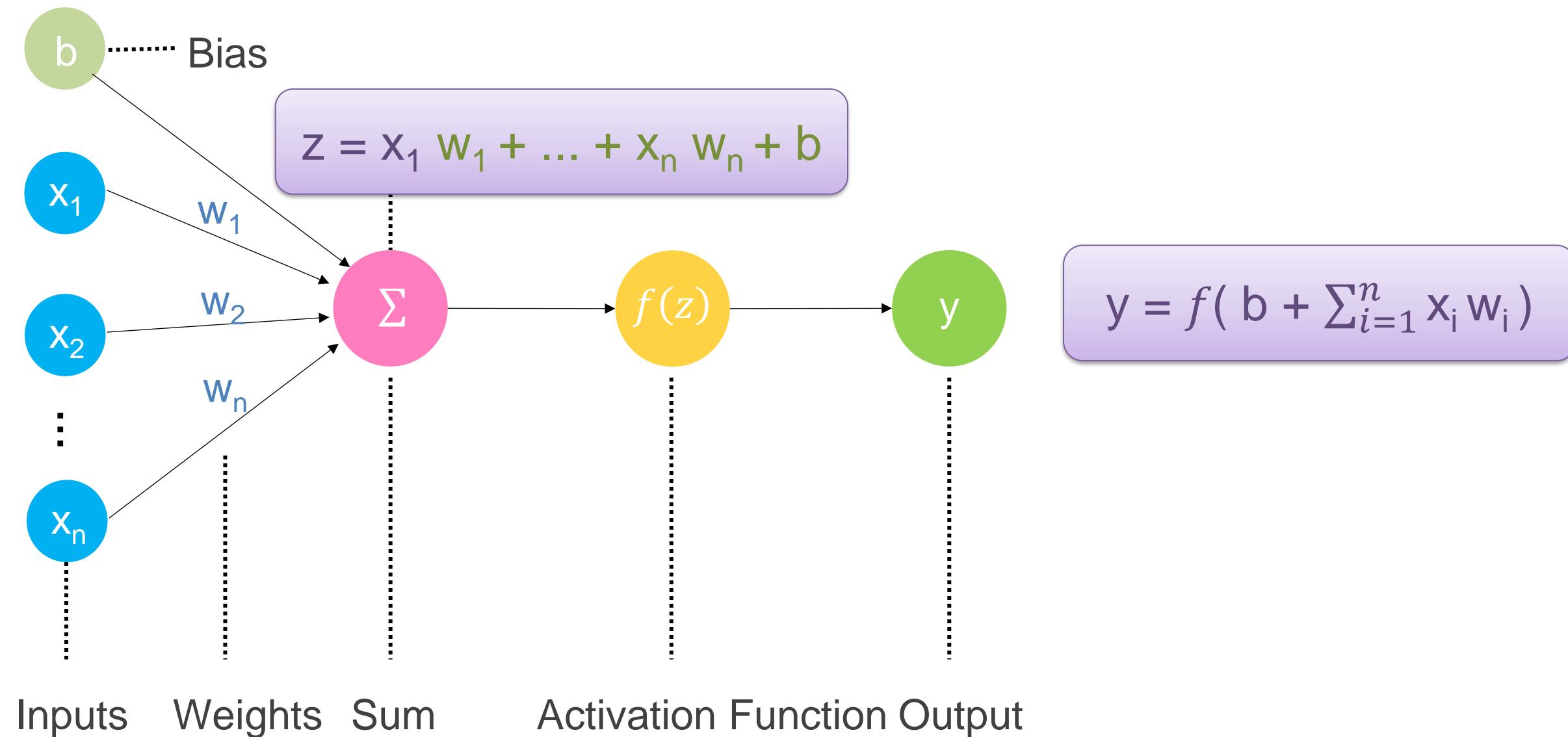
The aim is to minimize the cost to get the best performing neural network!

Improving Performance of a Neural Network



Forward Propagation

Forward Propagation



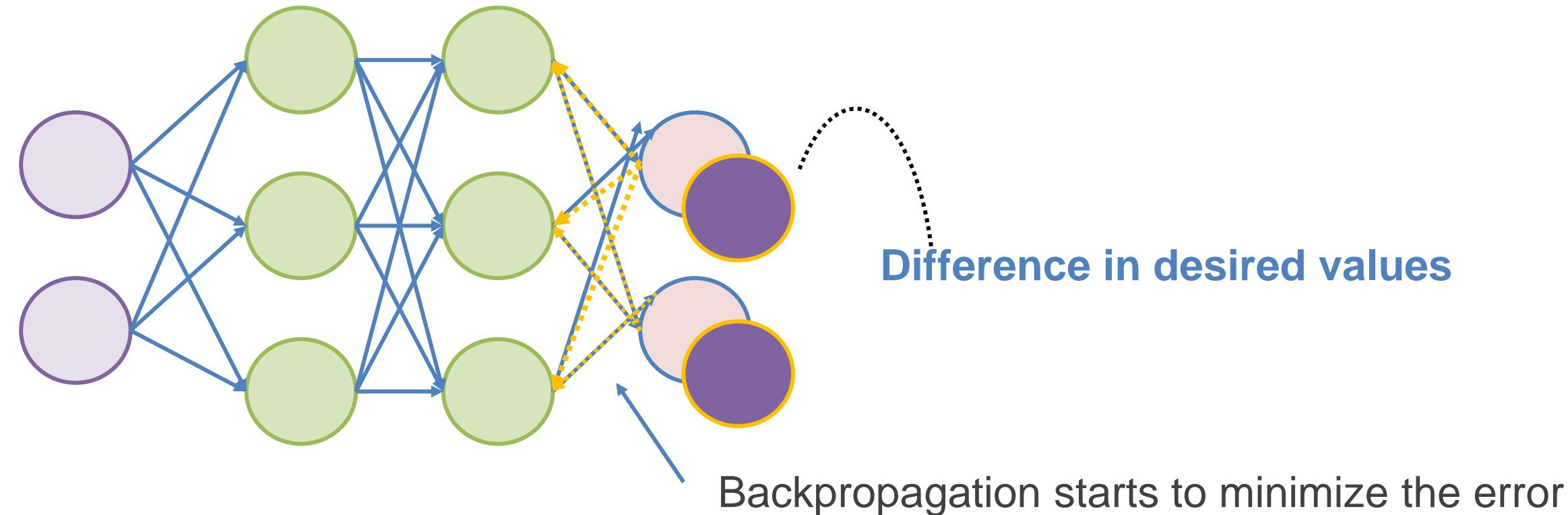
Understanding Backpropagation Using Gradient Descent

What Is Backpropagation?

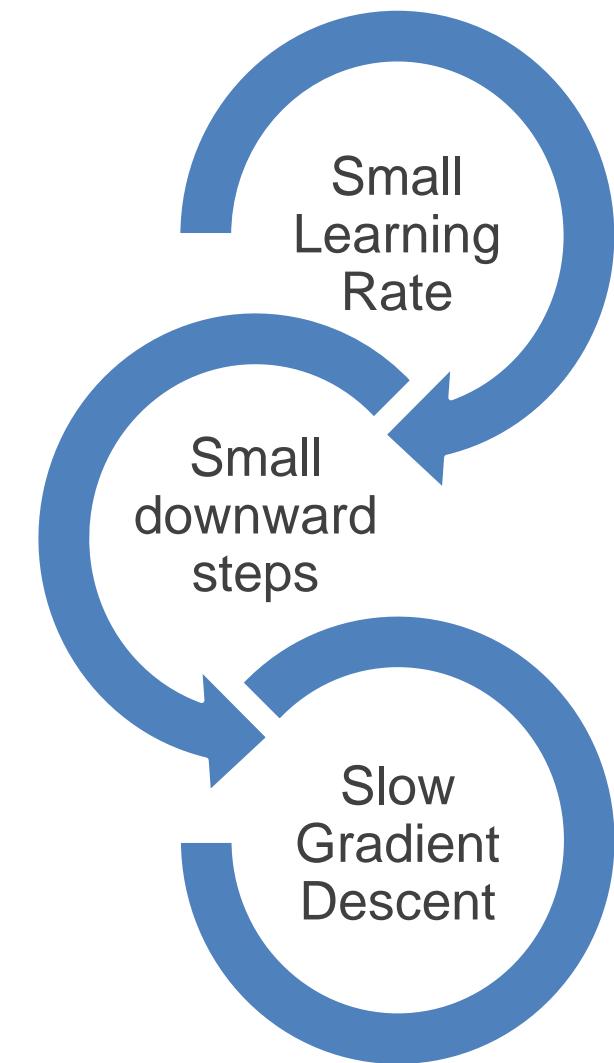
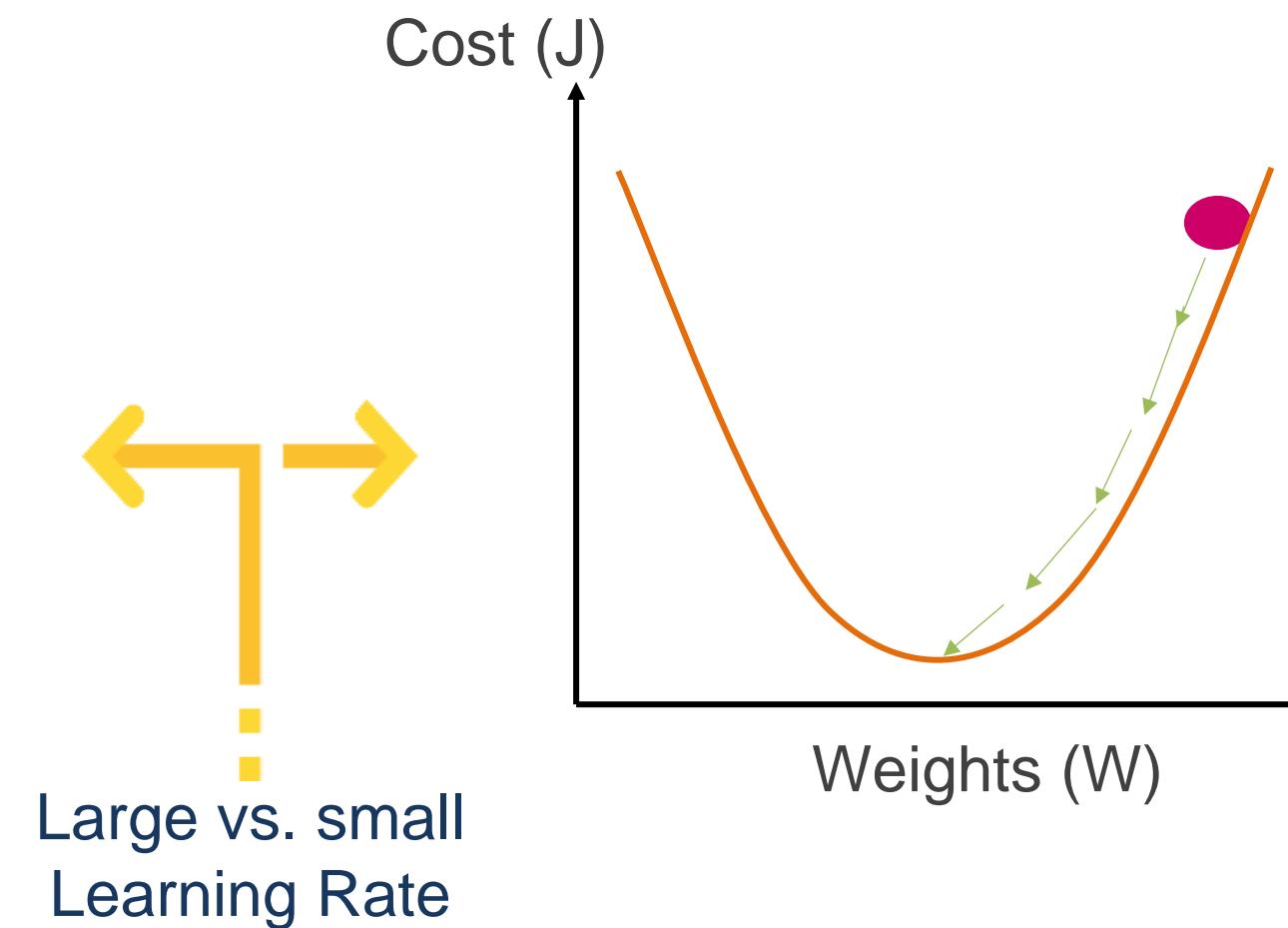
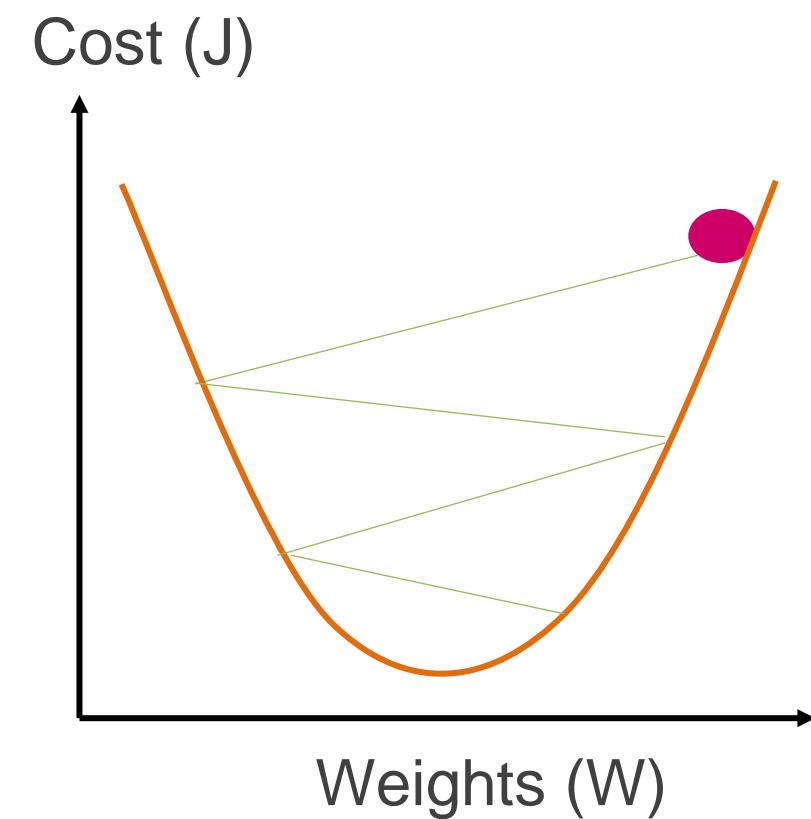
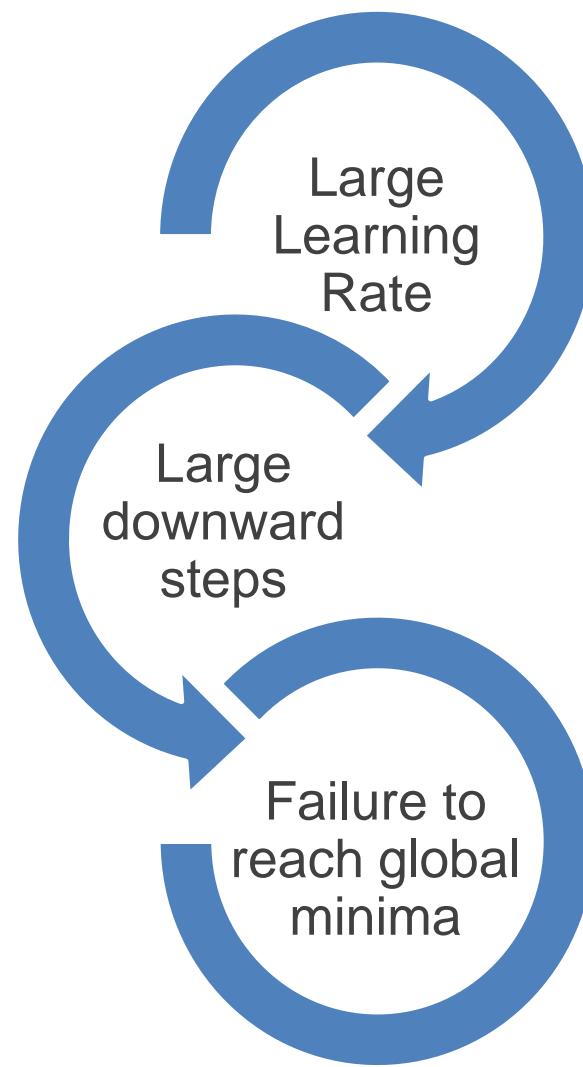
Algorithm to improve the accuracy of prediction and quick calculation of derivatives



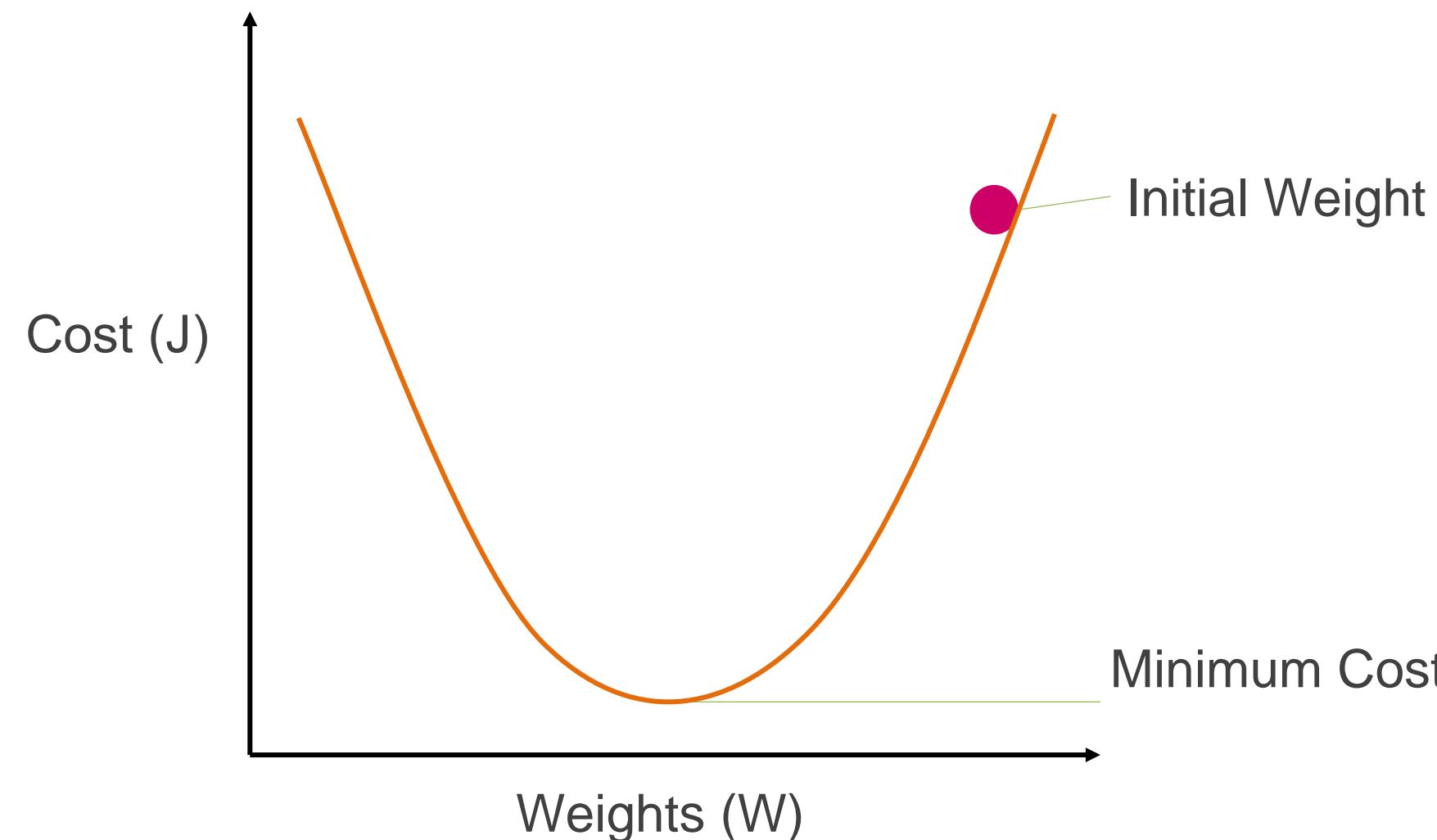
- e! A Learning algorithm to compute gradient descent with respect to weight
- e! Weights values are manipulated to get the predicted result
- e! Backpropagation is used to update the weights from output to input



Learning Rate (α)



Applying Gradient Descent To Cost Function



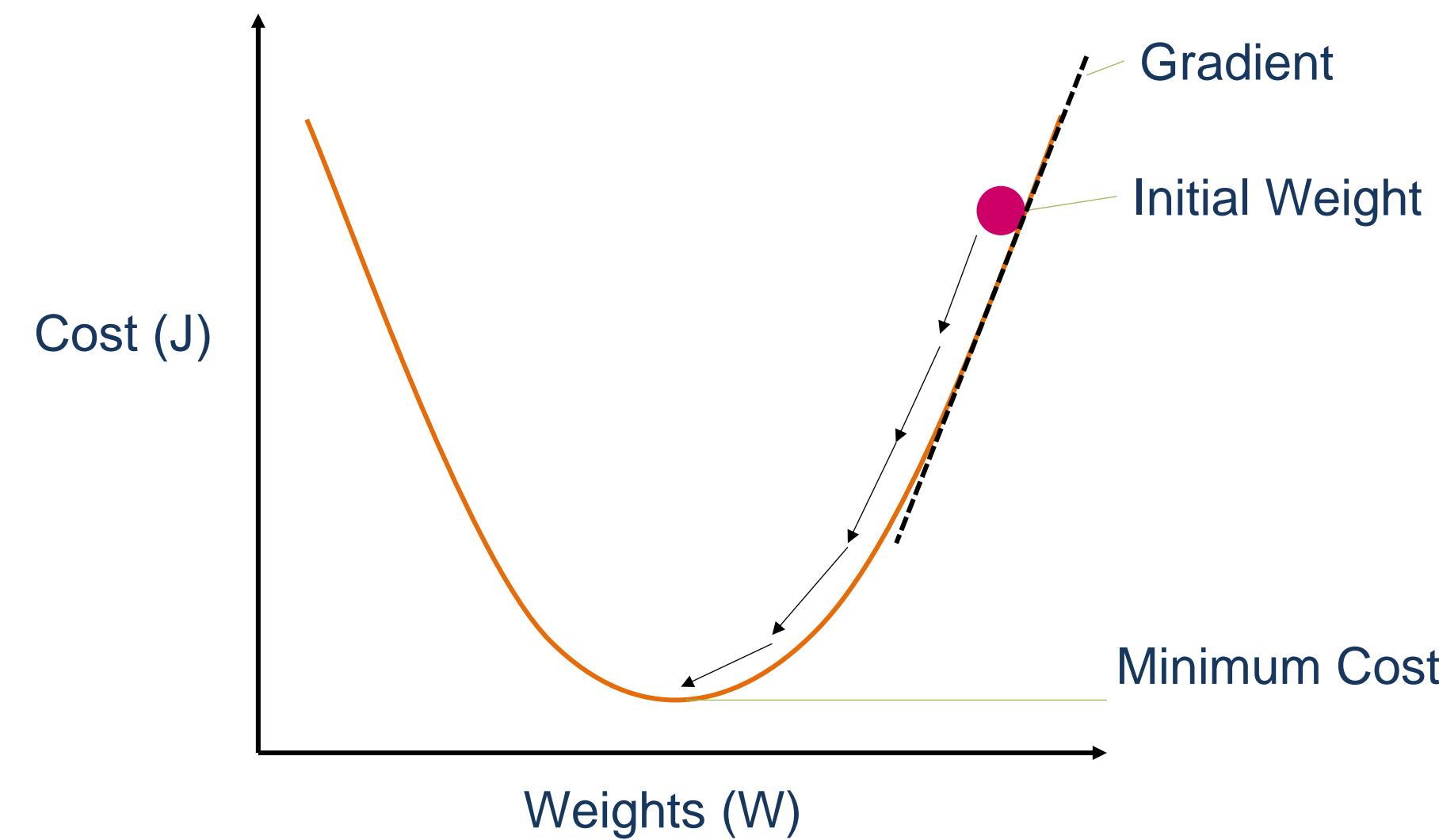
Gradients are used from moving from one point to another.

Thus, gradient of **Cost Function** is calculated with respect to **Initial Weight** to reach the

$$\text{Minimum Cost} - \frac{\partial J}{\partial W}$$

Gradient is simply differentiation of a function

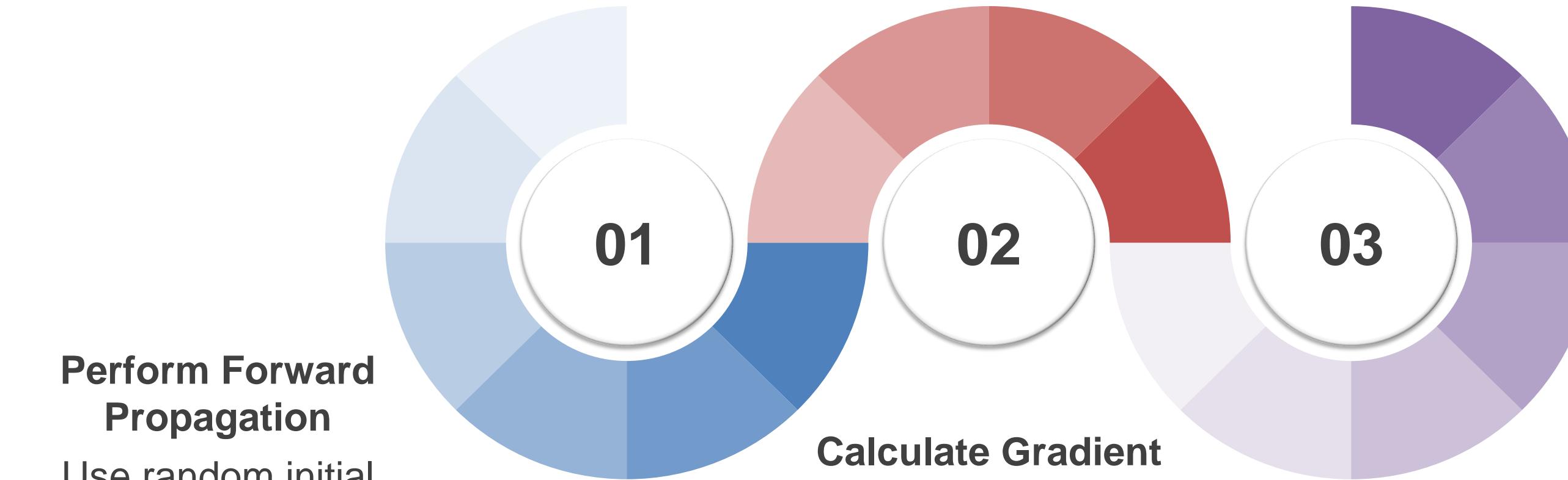
Applying Gradient Descent To Cost Function



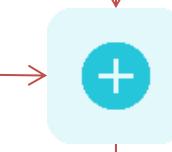
- e! Calculate gradient at each step and move in the direction of descend
- e! With each step, weights are updated to a value where the cost is minimum

Process Of Building A Neural Network

Process Of Building A Neural Network



$$W_{\text{new}} = W_{\text{old}} - \alpha \frac{\partial J}{\partial W}$$

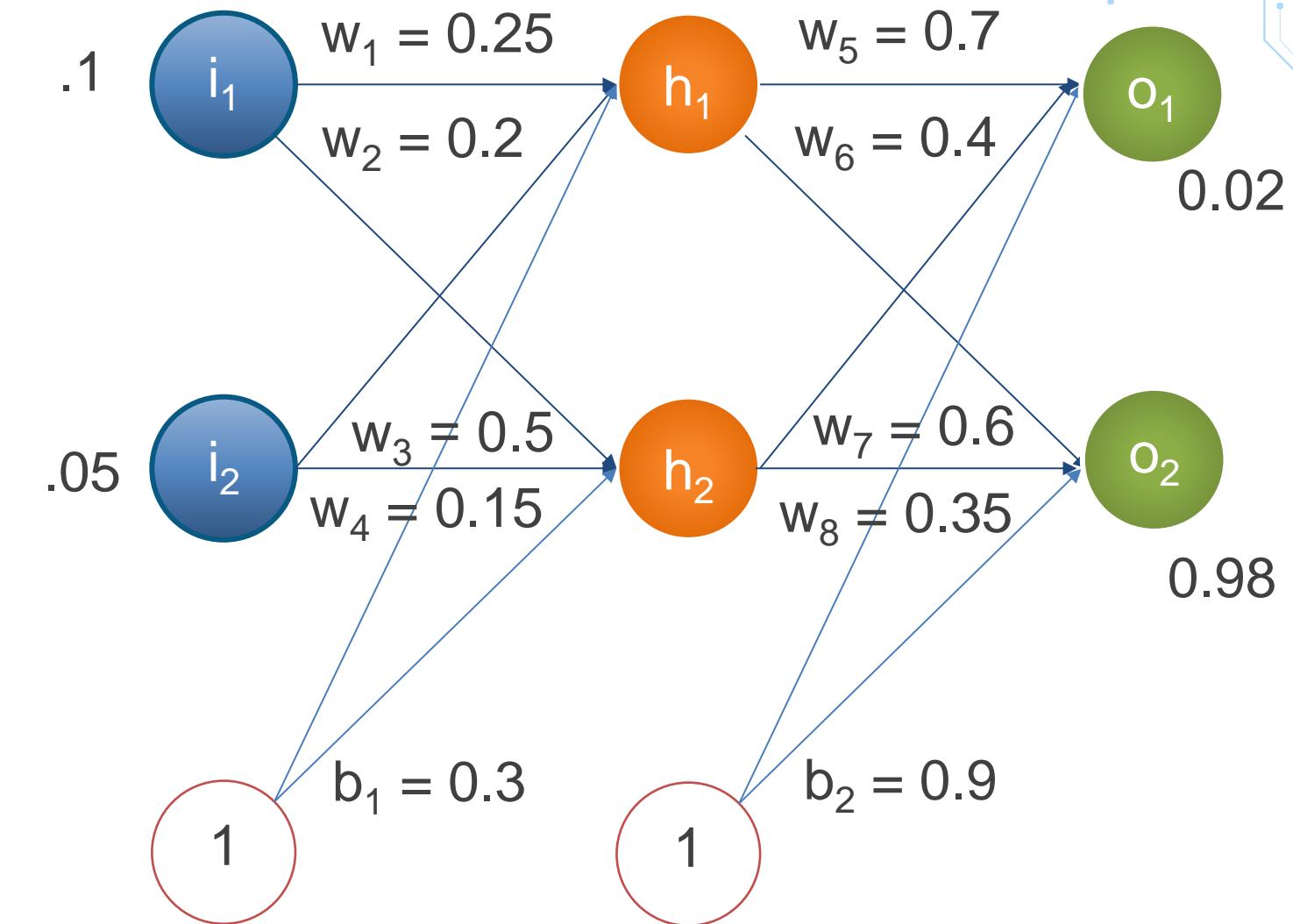


Backpropagation

Step-By-Step Learning Of A Neural Network

Consider the following neural network with:

- e! Two input neurons
- e! One hidden layer with two neurons
- e! Two output neurons



The goal of backpropagation is to optimize the weights so that the network can predict more accurate outputs for the given inputs

Learning Of A NN: Forward Propagation

e! Calculate the inputs to h1 and h2: Sum of weighted inputs

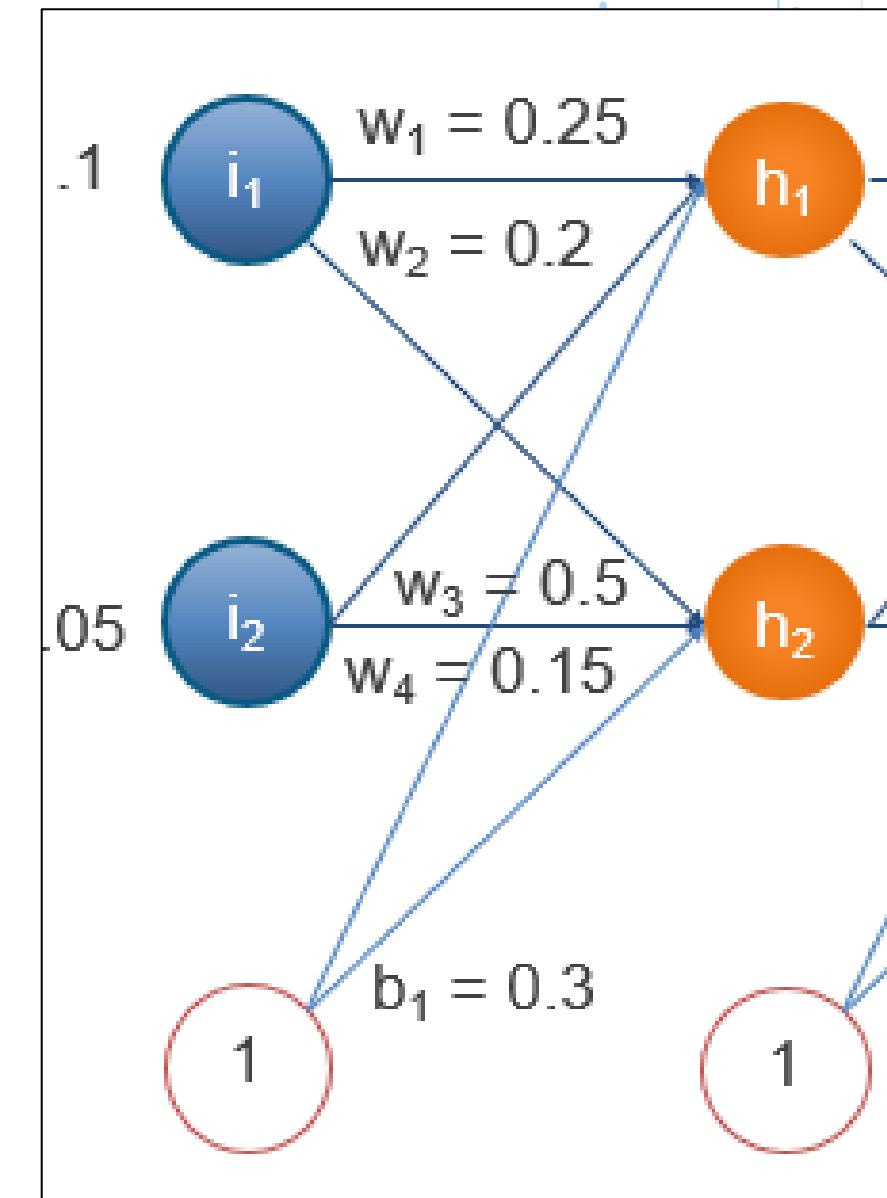
$$in_{h_1} = w_1 \cdot i_1 + w_3 \cdot i_2 + b_1 \cdot 1 = 0.25 \cdot 0.1 + 0.5 \cdot 0.05 + 0.3 = 0.35$$

$$in_{h_2} = w_2 \cdot i_1 + w_4 \cdot i_2 + b_1 \cdot 1 = 0.2 \cdot 0.1 + 0.15 \cdot 0.05 + 0.3 = 0.3275$$

e! Calculate the outputs from h1 and h2: Application of Activation Function

$$out_{h_1} = \frac{1}{1 + e^{-inh_1}} = 0.5866$$

$$out_{h_2} = \frac{1}{1 + e^{-inh_2}} = 0.5812$$



Learning Of A NN: Forward Propagation(Contd)

- e! Calculate the inputs to o1 and o2: Sum of weighted outputs from hidden layer

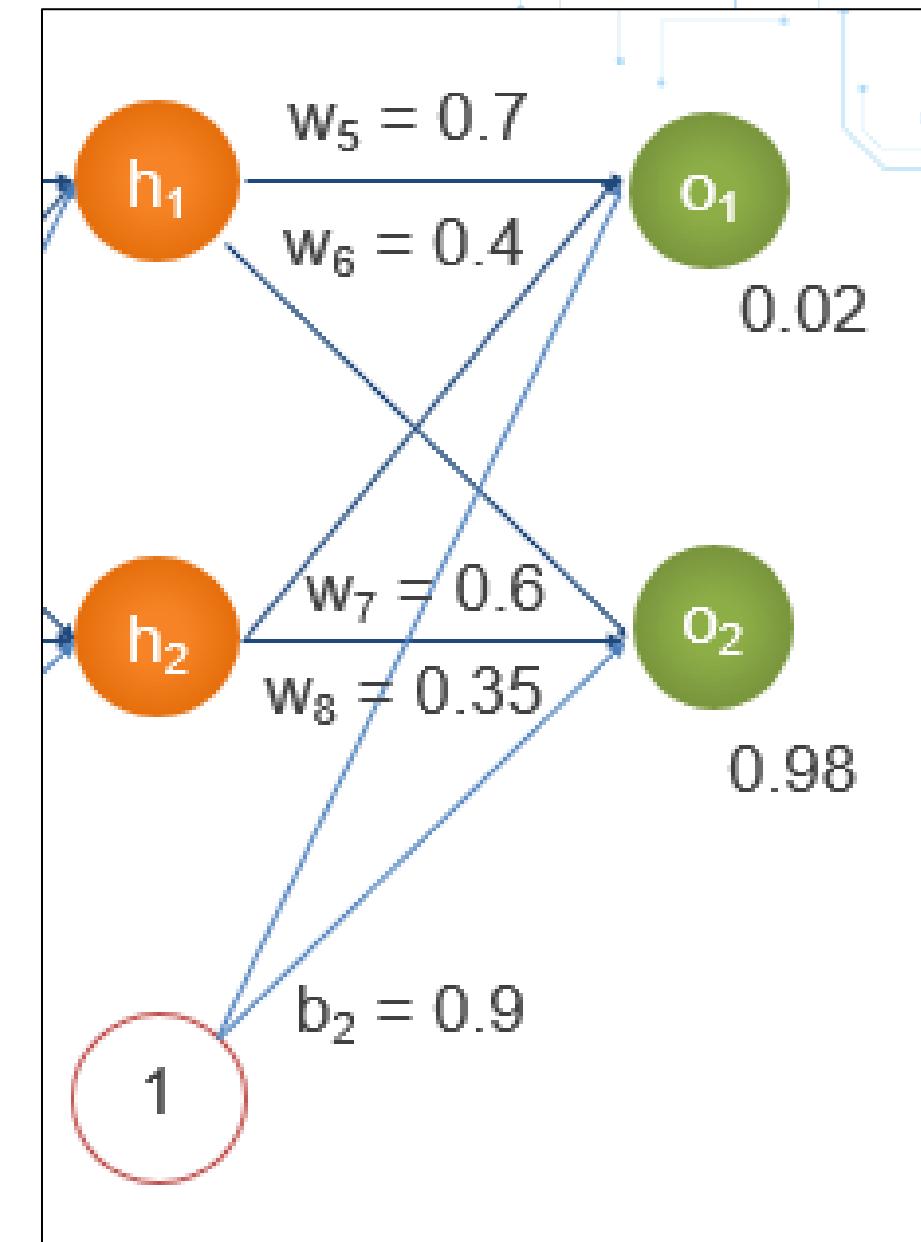
$$in_{o_1} = w_5 \cdot h_1 + w_7 \cdot h_2 + b_2 \cdot 1 = 1.3415$$

$$in_{o_2} = w_6 \cdot h_1 + w_8 \cdot h_2 + b_2 \cdot 1 = 1.1546$$

- e! Calculate the outputs from o1 and o2: Application of Activation function

$$out_{o_1} = \frac{1}{1+e^{-ino_1}} = 0.7927$$

$$out_{o_2} = \frac{1}{1+e^{-ino_2}} = 0.7603$$



Learning Of A NN: Calculating The Error

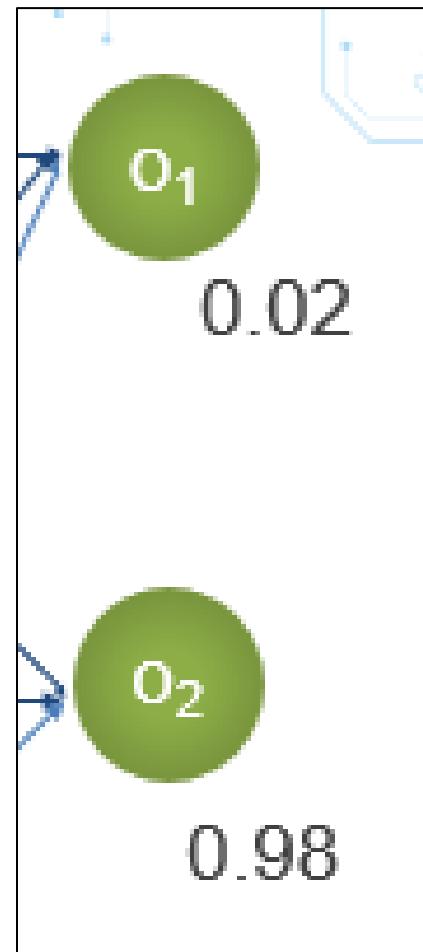
e! Calculate the errors for o1 and o2

$$E_{o_1} = \frac{1}{2}(\text{target}_{o_1} - \text{out}_{o_1})^2 = 0.2985$$

$$E_{o_2} = \frac{1}{2}(\text{target}_{o_2} - \text{out}_{o_2})^2 = 0.0241$$

e! Calculate the total error as shown below

$$E_{\text{total}} = E_{o_1} + E_{o_2} = 0.3226$$



Learning Of A NN: Backpropagation

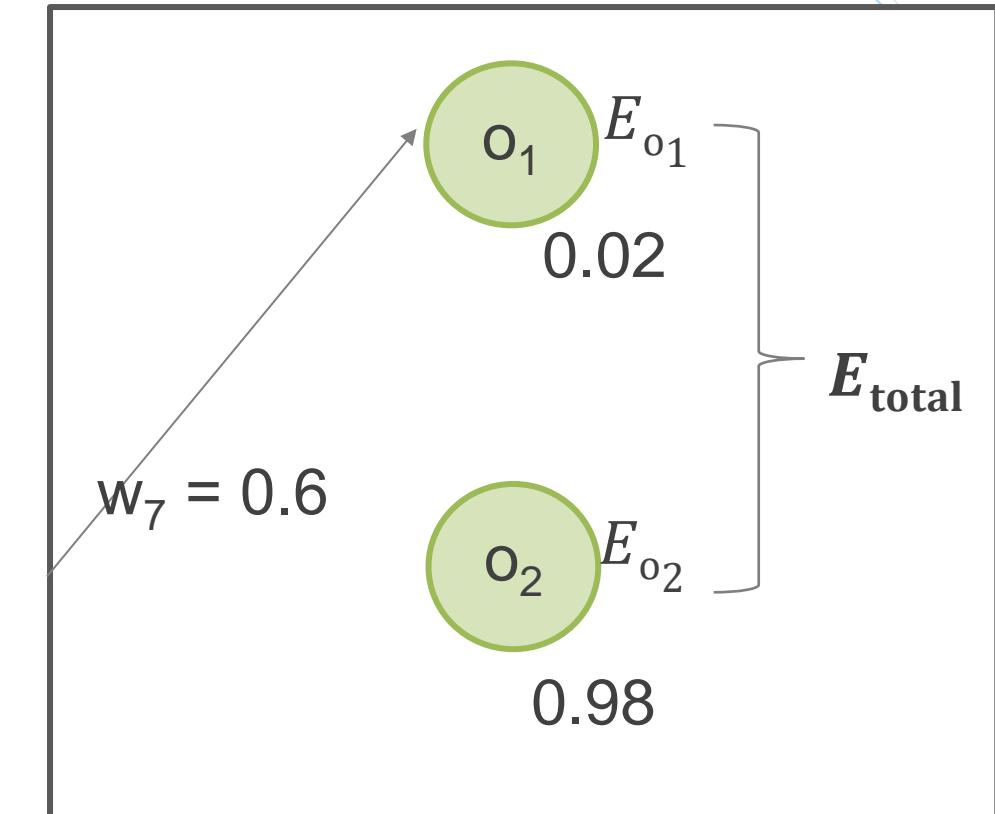
e! Calculate the error gradient of E_{total} with respect to w_7

e! Using Chain Rule, it can be calculated as:

$$\frac{\partial E_{total}}{\partial w_7} = \frac{\partial E_{total}}{\partial out_{o_1}} * \frac{\partial out_{o_1}}{\partial in_{o_1}} * \frac{\partial in_{o_1}}{\partial w_7} = 0.0745$$

e! w_7 is updated as follows:

$$w_7^{new} = w_7^{old} - \alpha * \frac{\partial E_{total}}{\partial w_7} = 0.5628 \dots \dots \dots \text{(Using } \alpha = 0.5\text{)}$$



Similarly, all other weights are updated

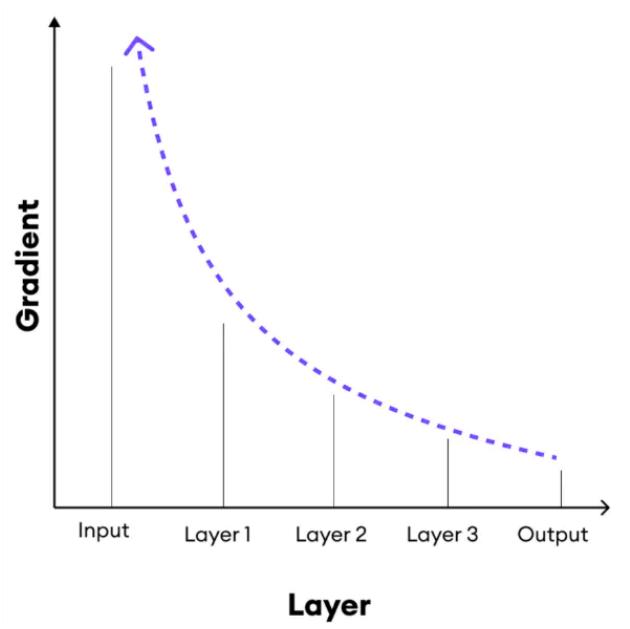
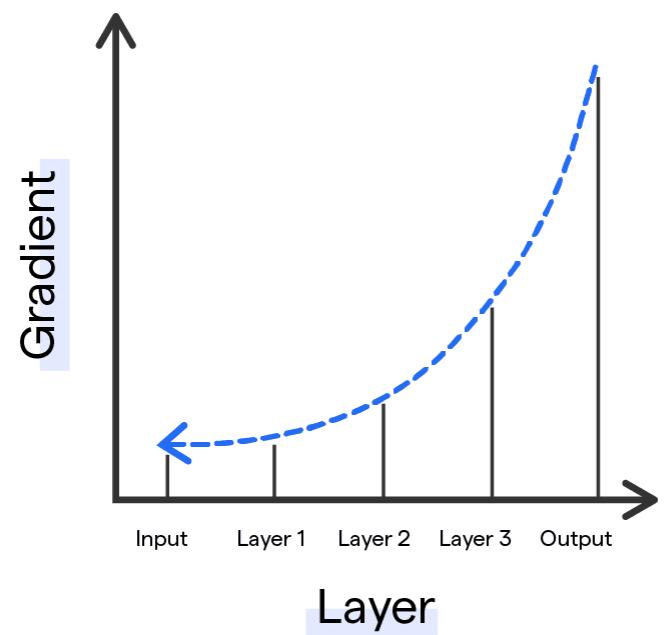
Weight Initialization Strategies

Weight initialization

What happens if there is no good initialization?



Slow convergence



Zero and Constant Initialization



Point	Zero Initialization	Constant Initialization
Definition	All weights are initialized to exactly zero .	All weights are initialized to a constant non-zero value
Technique	Set every weight $w = 0$.	Set every weight $w = c$ where c is a constant.
Purpose	Start with simplicity, but it fails for training .	Introduces small differences but still bad for deep learning .
Problem	No symmetry breaking → all neurons learn the same thing → no learning .	Same as zero: neurons behave identically → poor learning.
Best Used For	Not recommended for any serious models.	Only for experiments or extremely basic models.
Example	$W = 0$	$W = 0.01$ or $w = 0.5$

Random and Xavier/Glorot Initialization

Point	Random Initialization	Xavier or Glorot Initialization
Definition	Weights are initialized with random values	Weights are initialized randomly , but with a carefully calculated variance based on network size.
Technique	Randomly pick values without considering input/output sizes.	Randomly pick values with variance $\frac{2}{n_{in}+n_{out}}$, balancing inputs and outputs.
Purpose	Breaks symmetry better than zero but can still cause problems .	Maintain stable activations and gradients across layers for faster and safer training .
Problem	Can cause vanishing or exploding	Solves vanishing/exploding problems by scaling weights.
Best Used For	Very shallow or experimental networks.	Deep neural networks, with sigmoid or tanh activations.
Example	$w \sim U(-1, 1)$ or $w \sim N(0, 1)$	$w \sim U\left(-\sqrt{\frac{6}{n_{in}+n_{out}}}, \sqrt{\frac{6}{n_{in}+n_{out}}}\right)$

Shallow vs. Deep Neural Networks

Shallow vs. Deep Neural Networks

Aspect	Shallow Neural Networks	Deep Neural Networks
Definition	Neural networks with few layers	Neural networks with many layers
Complexity	Low complexity	High complexity
Learning Capacity	Limited learning capacity	Higher learning capacity
Risk of Overfitting	Lower	Higher
Data Requirement	Requires less data	Requires more data for effective training
Parameter Count	Fewer parameters	Many more parameters
Computational Resources	Requires less computation	Needs more resources (e.g., GPUs)
Interpretability	Easier to interpret	More difficult to interpret
Examples	Single-layer Perceptron, Logistic Regression	CNNs, RNNs

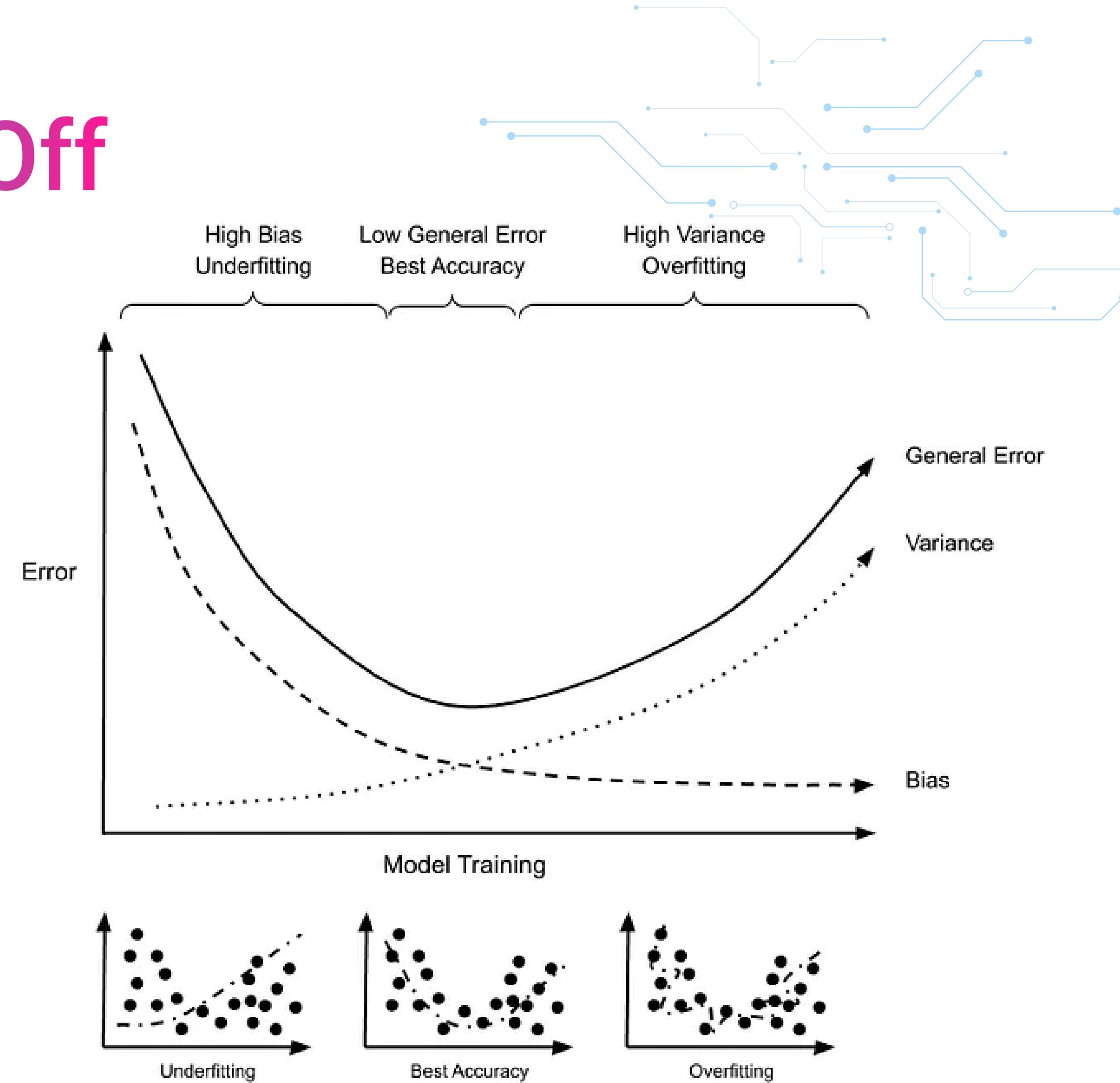
Bias-Variance Trade-off

Bias-Variance Trade Off

The bias-variance trade-off refers to the careful balance between **underfitting** and **overfitting**. The objective is to choose the **right model complexity** that enables it to perform well on new, unseen data.

In Deep Learning :

- e! Deep Neural Networks can easily **reduce bias** due to their complexity.
- e! But they are also **prone to high variance** if not regularized well.



Understanding Epochs, Batch Size and Learning Rate

Epochs, Batch Size and Learning Rate

Batch Size: The number of training samples processed together in one training step

Epoch: One full pass through the entire training dataset, composed of multiple iterations.

Learning Rate: A hyperparameter that determines model weights are adjusted during each update step.

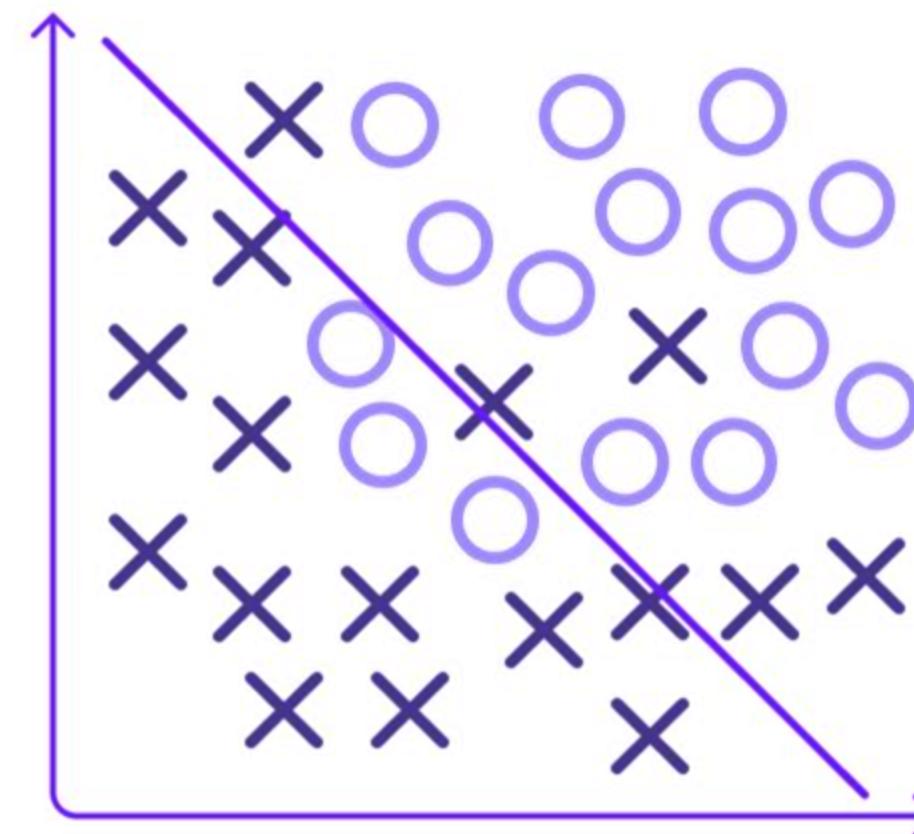
Let's consider a simple example using a feedforward neural network for image classification. Suppose we have a **dataset of 10,000 images**, and we want to train our model with the following parameters:

- e! **Batch size:** 1000
- e! **Number of epochs:** 50
- e! **Learning rate:** 0.01

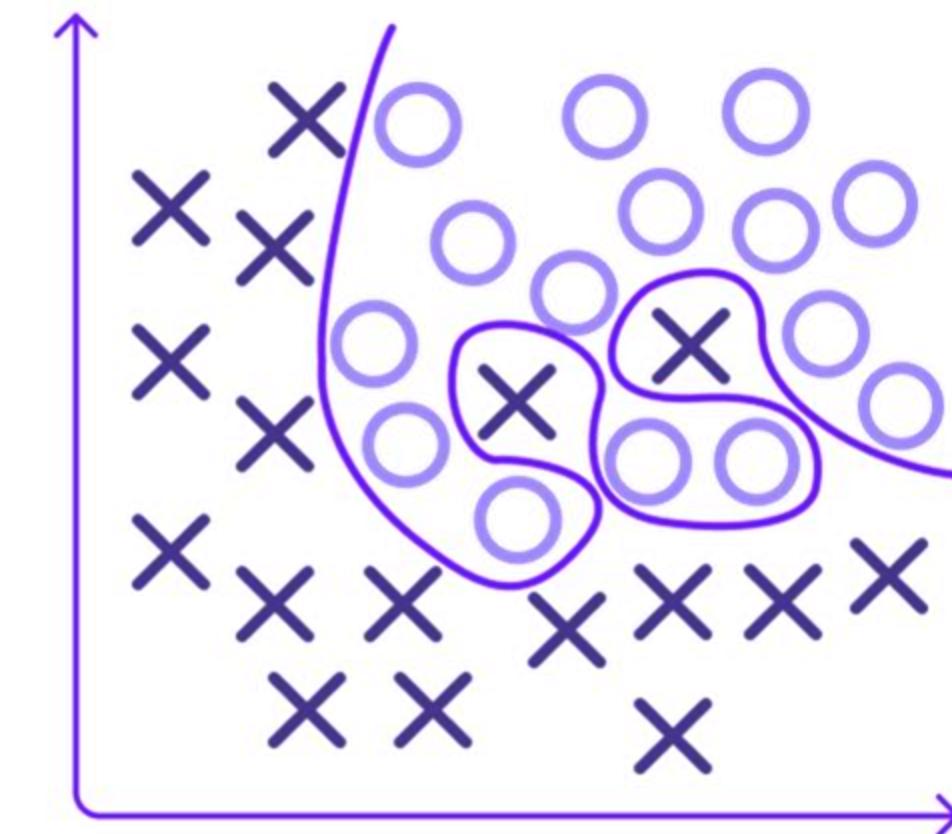
Overfitting and Underfitting

Overfitting and underfitting in deep learning

Underfitting occurs when a deep learning model is **too simple to capture** the underlying patterns in the training data.



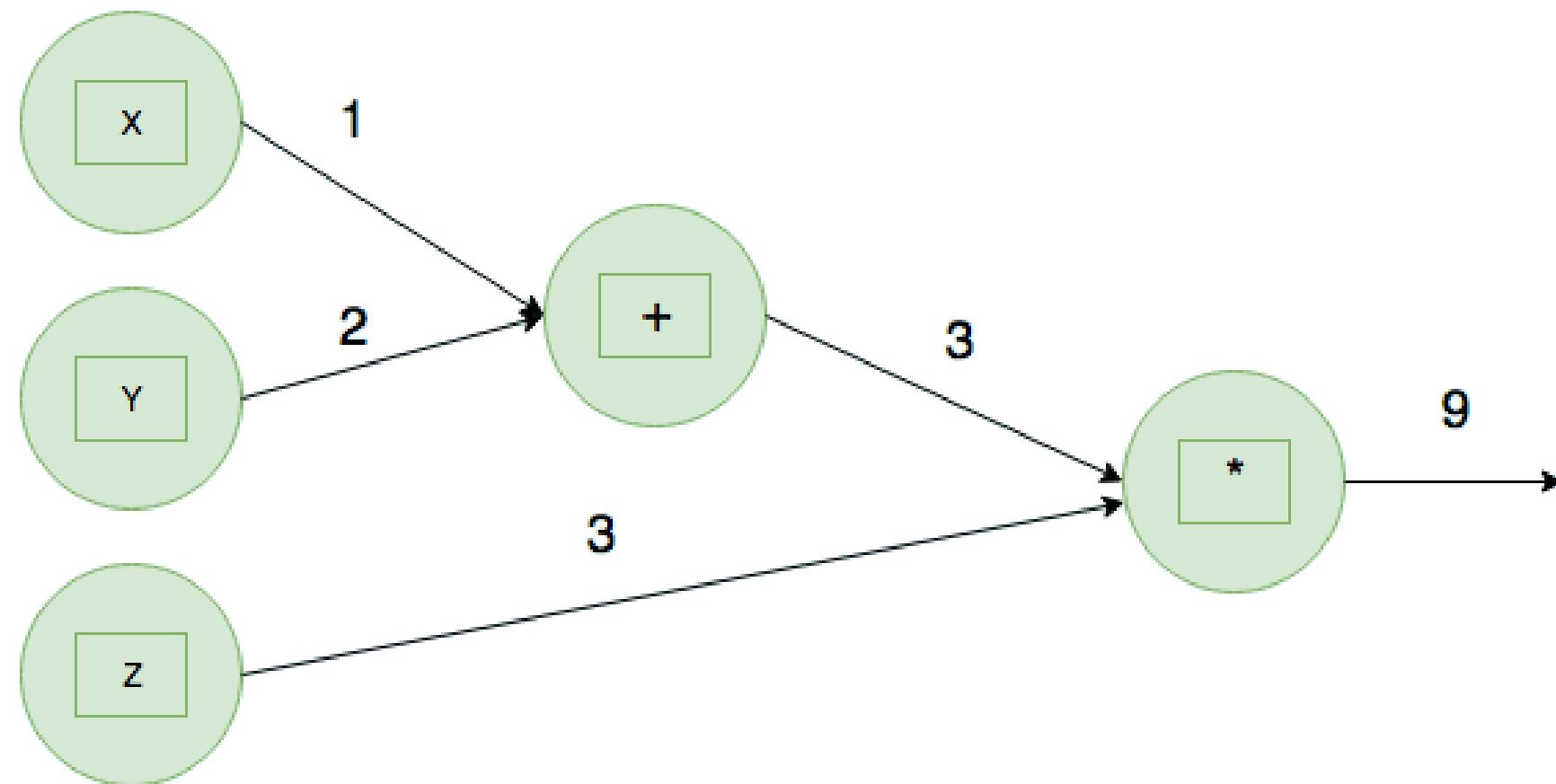
Overfitting happens when a deep learning model learns the **training data too well**, including **noise and irrelevant patterns**.



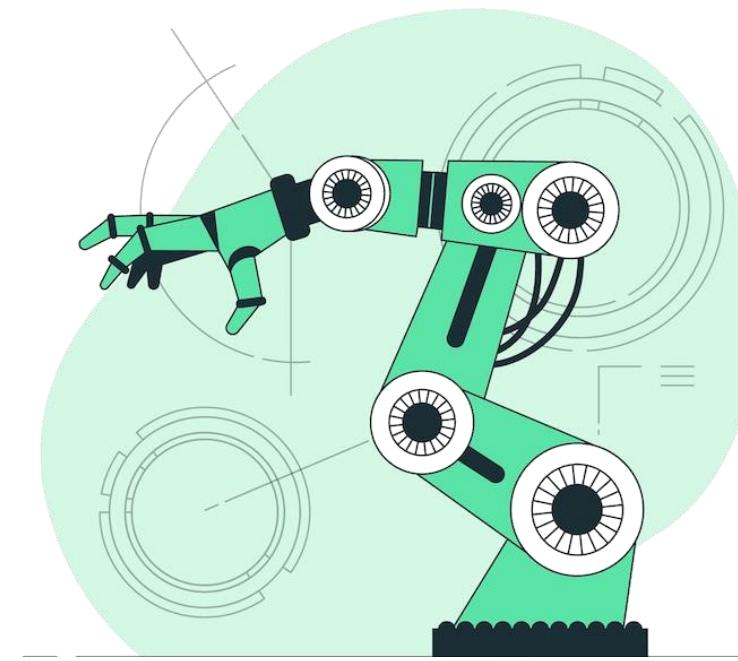
Computational Graphs in Deep Learning

Computational Graphs in Deep Learning

A computational graph in deep learning is a **structured representation** of the sequence of operations involved in a model, where nodes **represent mathematical operations** and edges represent the flow of data (tensors).



Why are they important?



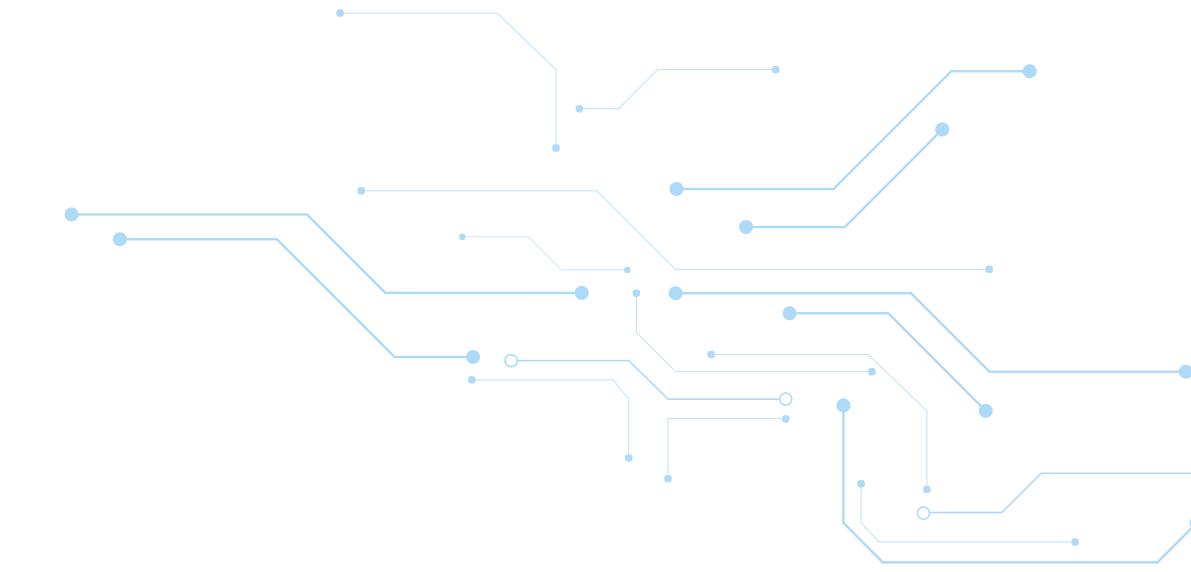
Automatic Differentiation



Modular Design



Optimization & Efficiency



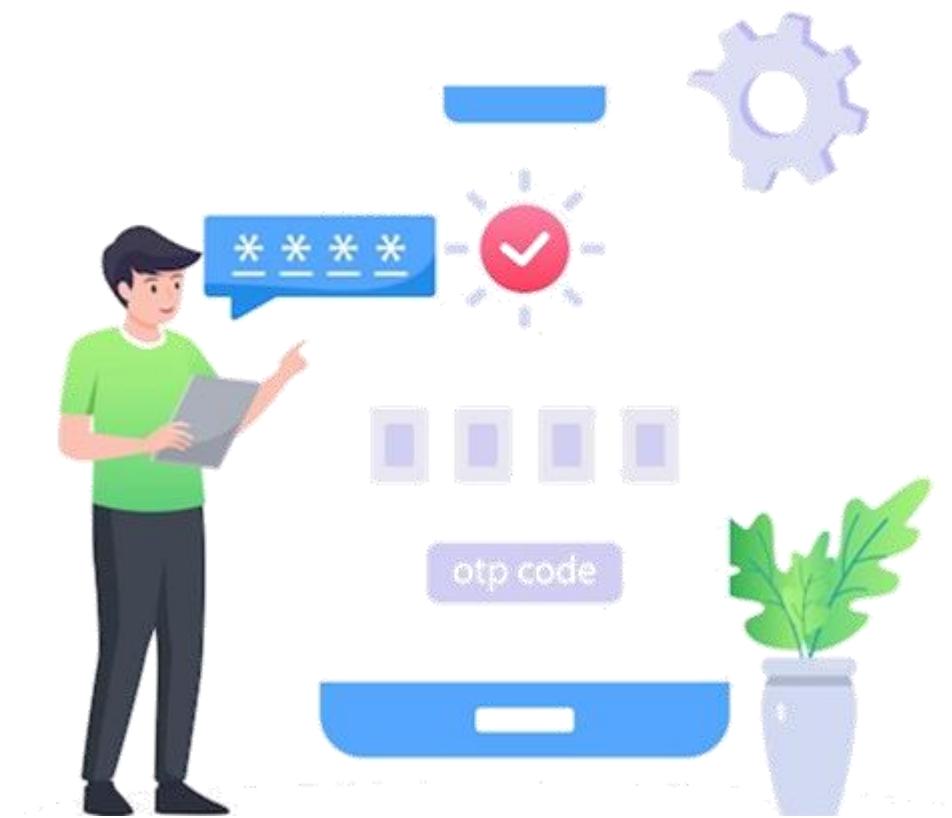
Dynamic & Static Graph Support

Monitoring Training and Validation Loss

Training and Validation Loss in Deep Learning

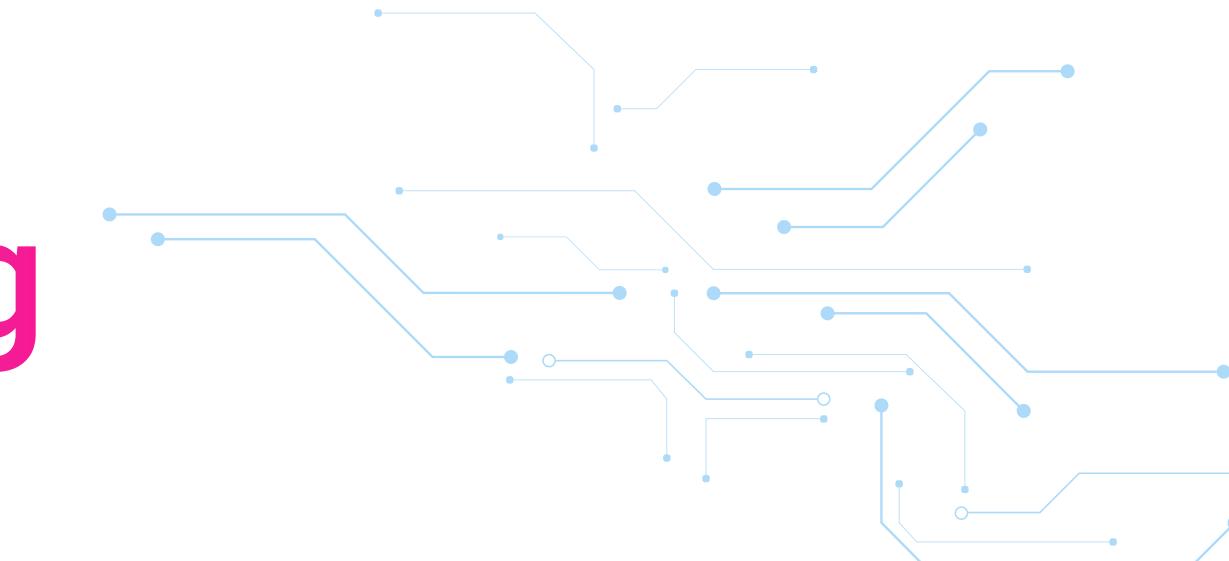


Training Loss measures the error on the training dataset during learning.



Validation Loss measures the error on a separate validation dataset (unseen during training).

Ideal Behavior While Monitoring



Situation	Training Loss	Validation Loss	Interpretation
Good Fit	Decreasing	Decreasing	Model is learning correctly
Overfitting	Decreasing	Increasing	Model is memorizing, not generalizing
Underfitting	High	High	Model is too simple or not trained enough

Creating a simple DNN using Keras

Simple DNN using Keras

Code example

```
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense
from sklearn.datasets import make_classification

# Step 1: Generate simple binary classification data
X, y = make_classification(n_samples=100, n_features=4, n_classes=2, random_state=42)

# Step 2: Define the simplest possible DNN
model = Sequential()
model.add(Dense(4, activation='relu', input_shape=(4,))) # Hidden layer
model.add(Dense(1, activation='sigmoid')) # Output layer

# Step 3: Compile and train the model
model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])
model.fit(X, y, epochs=2, batch_size=10)
```

Output

```
Epoch 1/2
10/10 ━━━━━━━━━━━━ 4s 14ms/step - accuracy: 0.5475 - loss: 0.6564
Epoch 2/2
10/10 ━━━━━━━━━━━━ 0s 10ms/step - accuracy: 0.6058 - loss: 0.6476
```

Implementing a perceptron model (Demonstration)

Note: Refer to the Module 1: Demo 1 on LMS for detailed steps.

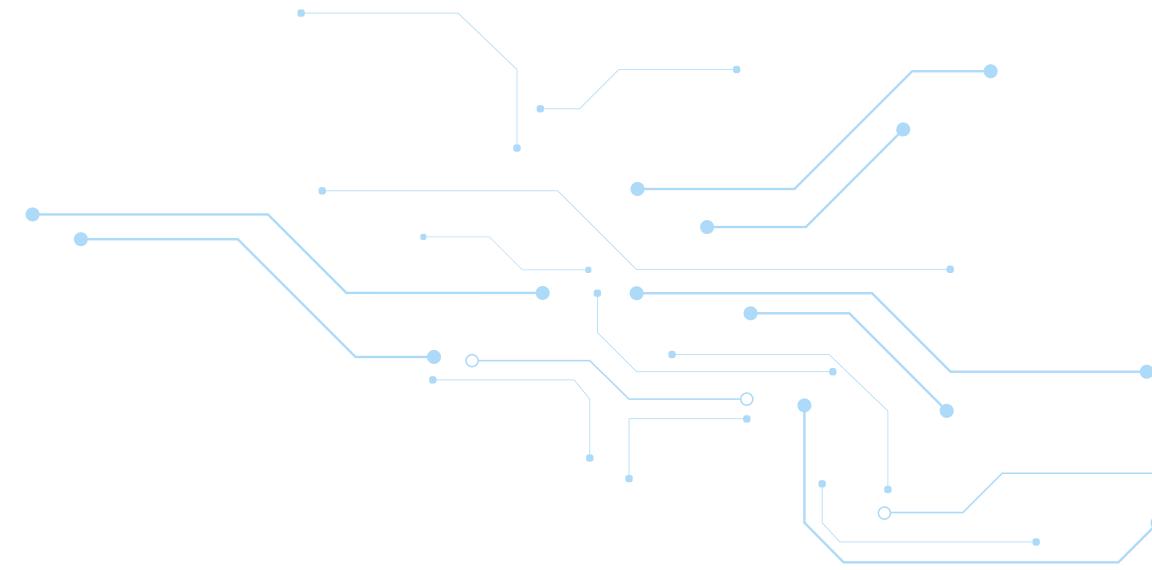
Training a basic MLP on the MNIST dataset (Demonstration)

Note: Refer to the Module 1: Demo 2 on LMS for detailed steps.

Summary

In this lesson, you have learned to:

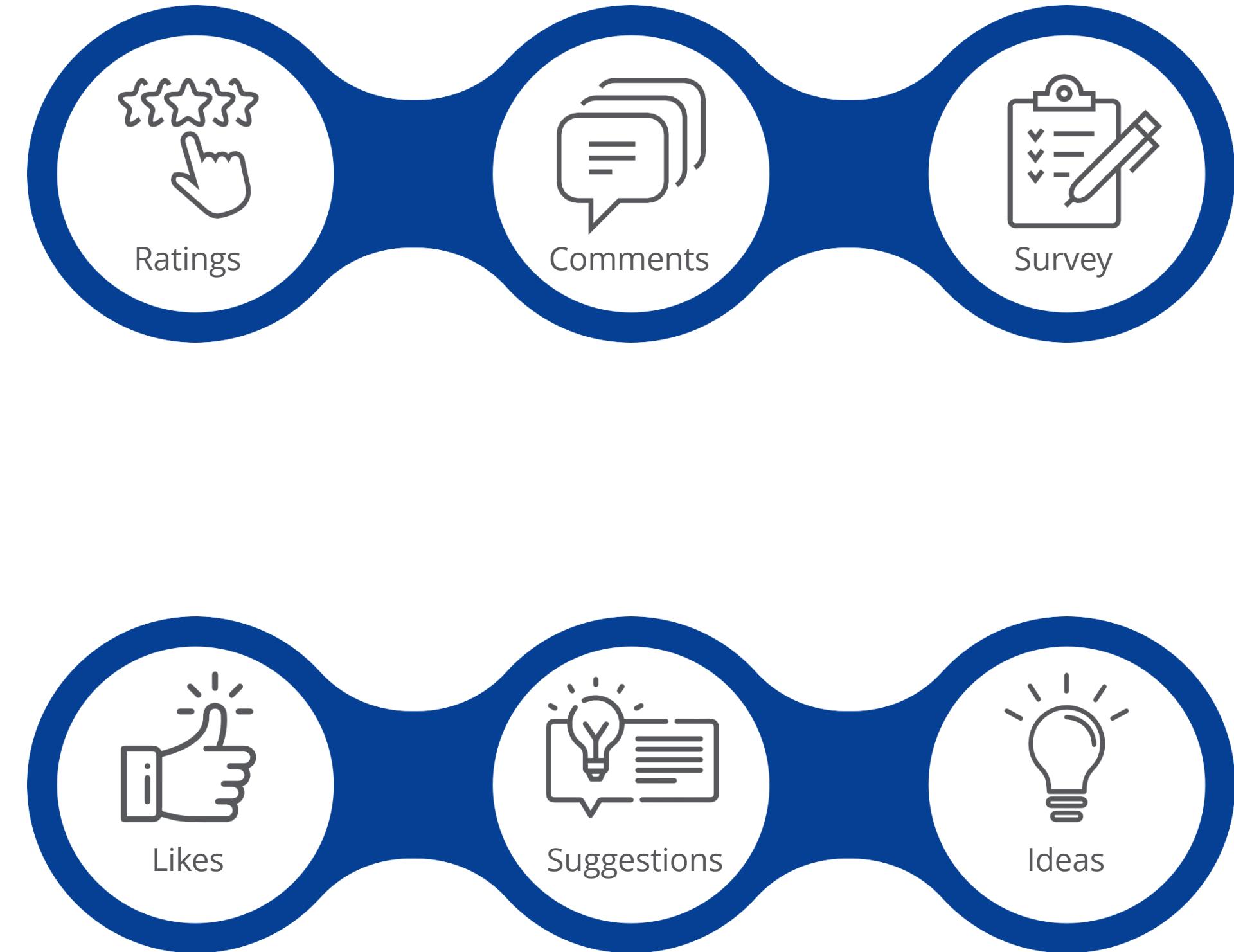
- e! Explain NN basics: Perceptron, MLP, activation & loss functions.
- e! Understand forward/backpropagation & gradient descent.
- e! Grasp bias-variance, overfitting, and key training terms.
- e! Use Keras to build and monitor a simple DNN..



Questions



Feedback



Thank You

For information, Please Visit our Website
www.edureka.co

