# Outline

- Introduction
- Learning paradigms
- History of artificial neural networks (ANN)
- Modelling of ANNs
- Multilayer perceptron (MLP)
- Gradient Descent and Backpropagation
- ANN types, design and issues
- Validation techniques for efficient learning
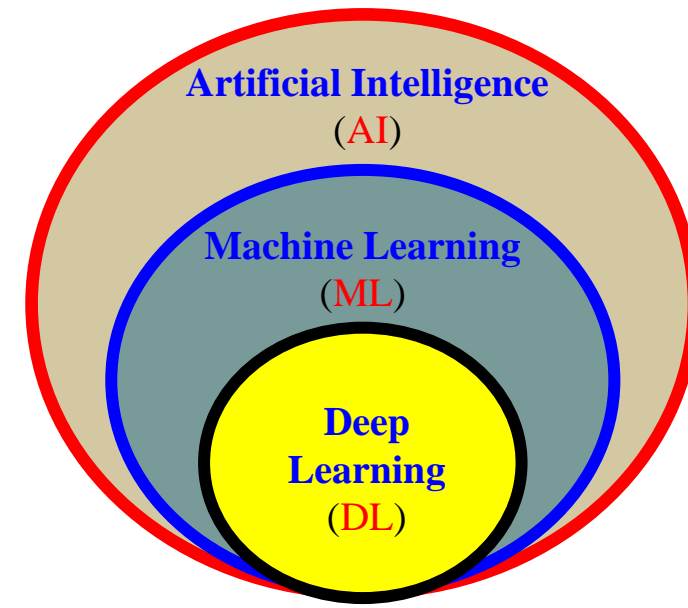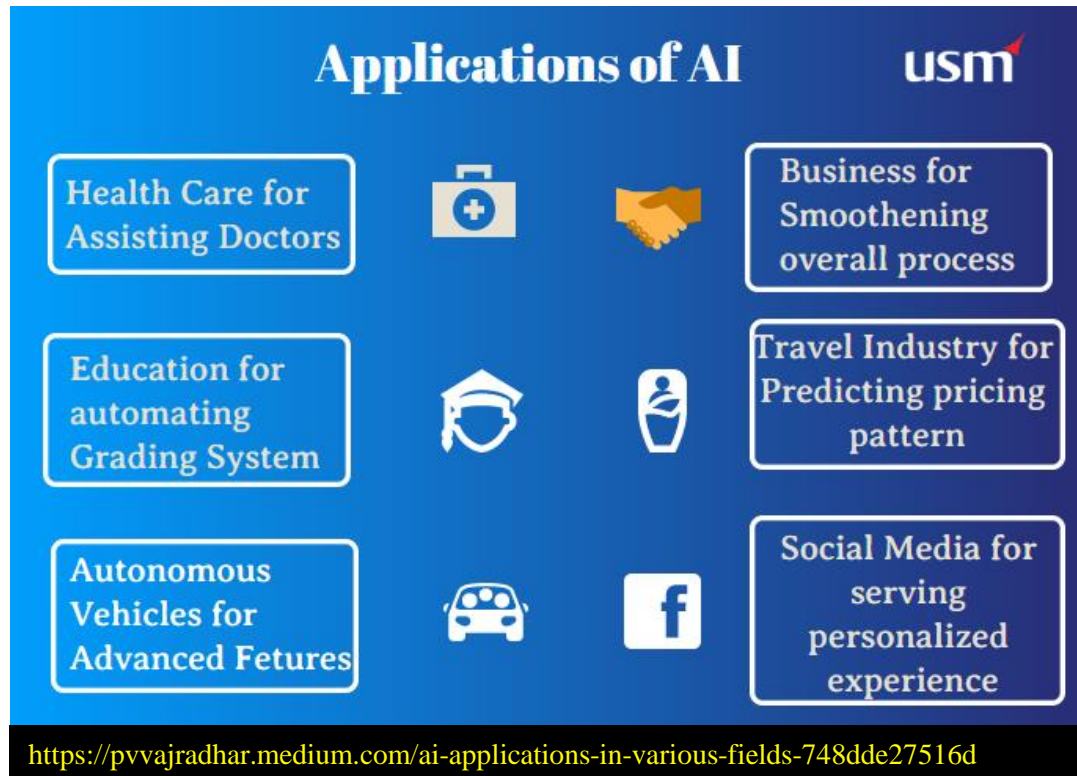- Assignment(s)
- Conclusion

# Introduction

❑ The ever-increasing popularity of artificial intelligence (AI) and machine learning (ML) provides a groundbreaking impetus on many aspects of our life.

➢ **Artificial Intelligence** (AI) are those set of human-designed tools (programs) to do things that is typically done by human

➢ **Machine learning** (ML) is an AI field where *machine* can learn new things *through* **experience** without the involvement of a human.

➢ **Deep learning** (DL) is a ML subset where machines adapt and learn from vast amount of data

# Categories of Machine Learning

## Learning Paradigms

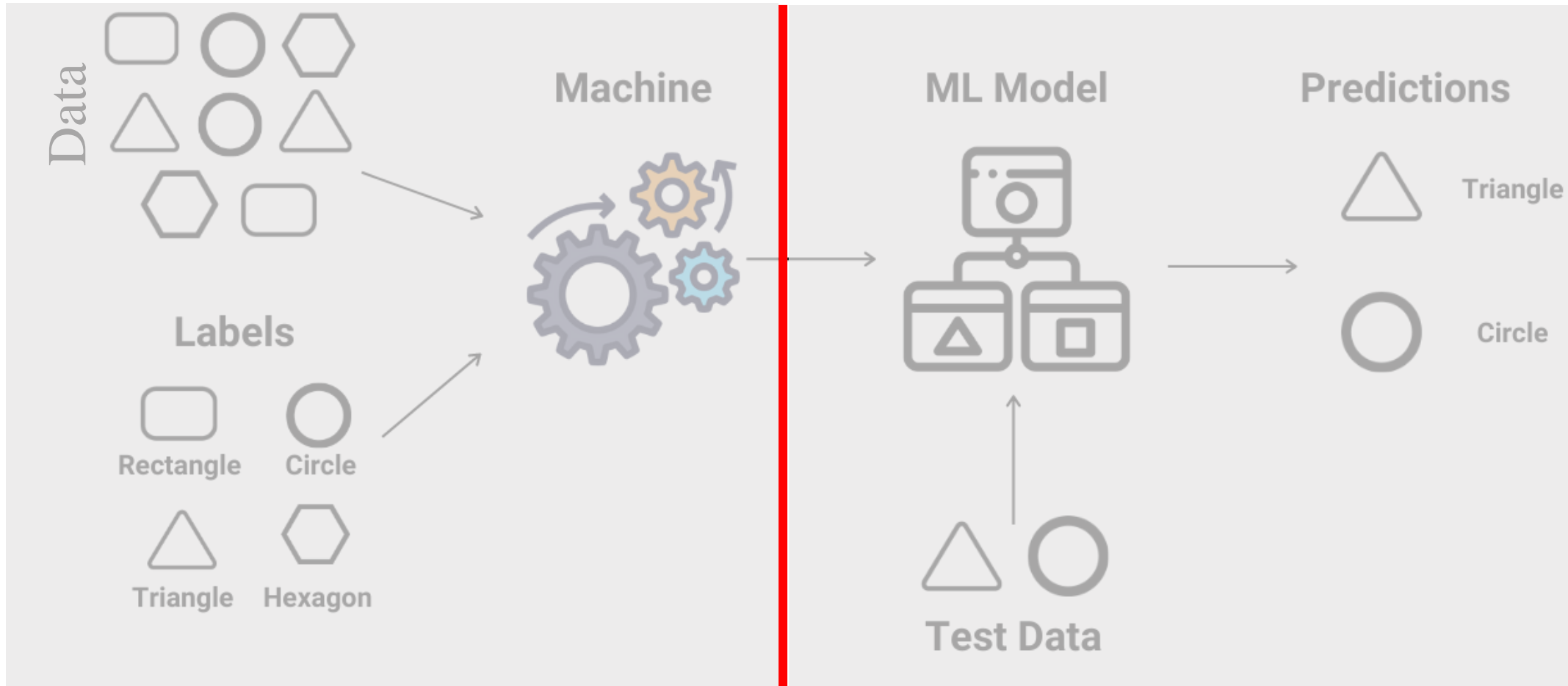| Supervised Learning | Reinforcement Learning | Unsupervised Learning |
|---|---|---|
| ➤ Learning with a teacher | ➤ Interactive learning environment by trial and error using feedback from its own actions and experiences. | ➤ Learning without a teacher |
| ➤ Data with known output (label) is given | | ➤ No labels |
| ➤ Classification and Regression | | ➤ Machine understand the data |
| | | ➤ Clustering |

| | | |
|---|---|---|
| Support Vector Machine (SVM), K Nearest Neighbours (KNN), Decision Trees, Random Forest *Feedforward* Artificial Neural Network (ANN) | Q-learning, Markov Decision Process | Gaussian Mixtures, *K*-means, RNN, Fuzzy *c*-means |

# Supervised Machine Learning

**Image source**: https://www.enjoyalgorithms.com/blog/classification-of-machine-learning-models

# Supervised Machine Learning (cont'd)



Training       Testing (or *validation*)

# Supervised Machine Learning (cont'd)

## Learning Paradigms

| Supervised Learning | Reinforcement Learning | Unsupervised Learning |
|---|---|---|
| ➢ Learning with a teacher<br>➢ Data with known output (label) is given<br>➢ Classification and Regression | ➢ Interactive learning environment by trial and error using feedback from its own actions and experiences. | ➢ Learning without a teacher<br>➢ No labels<br>➢ Machine understand the data<br>➢ Clustering |
| Support Vector Machine (SVM), K Nearest Neighbours (KNN), Decision Trees, Random Forest<br>*Feedforward* Artificial Neural Network (ANN) | Q-learning, Markov Decision Process | Gaussian Mixtures, *K*-means, RNN, Fuzzy *c*-means |

7

# History of Neural Networks (NN)

- **1940**: McCulloch and Pitts: *First* mathematical model of a *neuron* (A verification model)

- **1957**: Rosenblatt's: The *Perceptron* model

- **1959**: Widrow and Hoff developed MADALINE was the first NN to be applied to a *real-world* problem

**Progress on NN research halted until 1981**

- **1982**: Hopfield: Associative memory - Recurrent NN (or the RNNs)

- **1986**: Rumelhart: Backpropagation and the *era* of multilayer perceptron (MLP).

- **1990s**: Rise of support vector machine (SVM)

- **1997**: Schmidhuber & Hochreiter: An RNN, long short-term memory (LSTM) was proposed.

- **2006**: Hinton et al.: NN returned to the public's vision again though Deep belief nets (DBNs)

- **2016**: Boom of NN (Deep convolutional neural networks (CNNs): AlexNet, GoogLeNet, VGG, ResNet, etc.
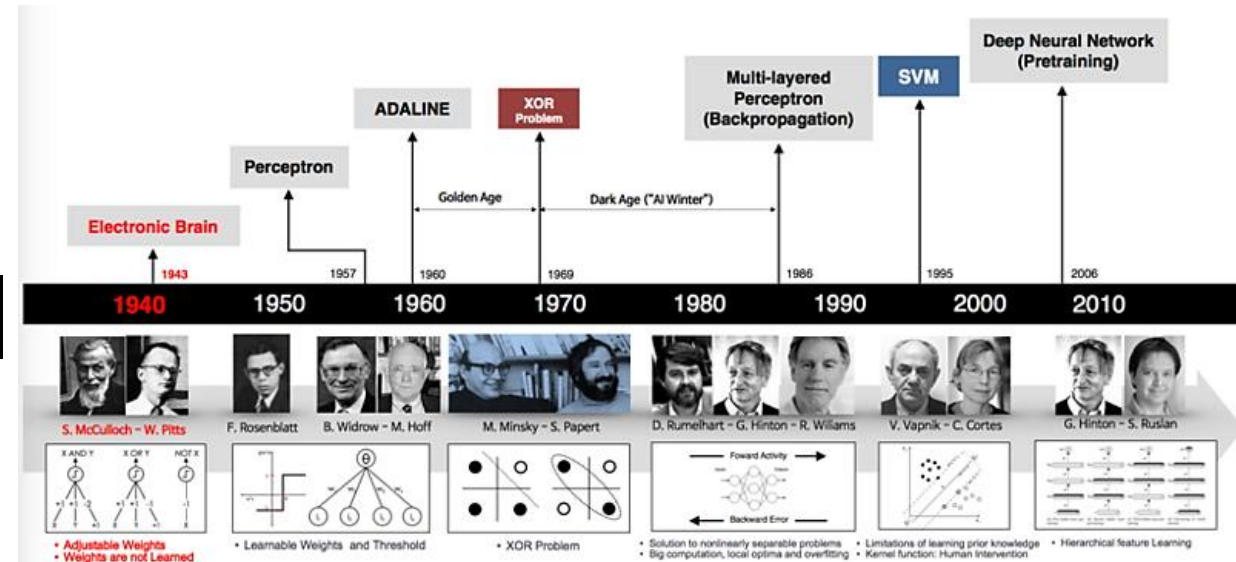


**Image source:** https://developpaper.com/take-you-into-the-past-life-and-this-life-of-neural-network/

8

# Human Brain and Biological Neurons

❏ Human brain contains billion of neurons (~10 billion)

❏ Each neuron is a cell that uses biochemical reactions to **receive**, **process** and **transmit** information

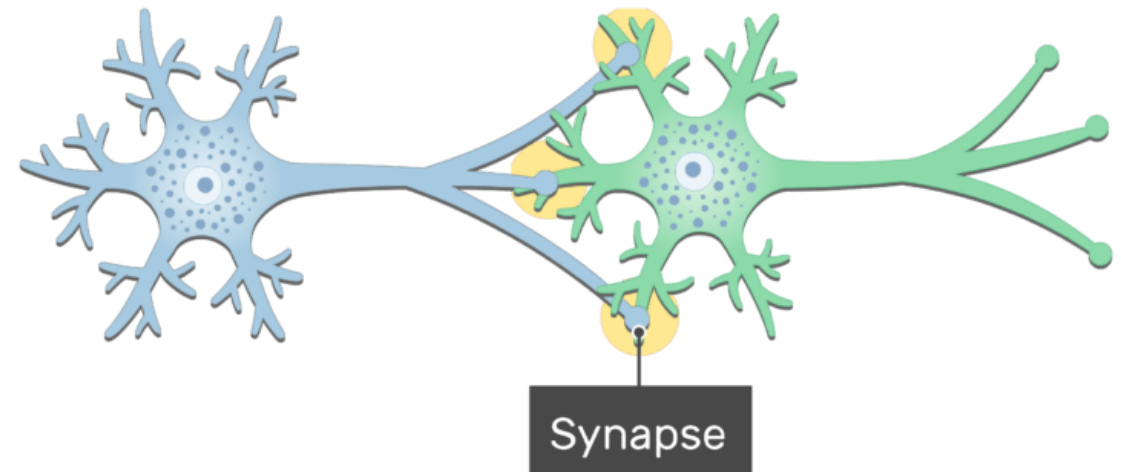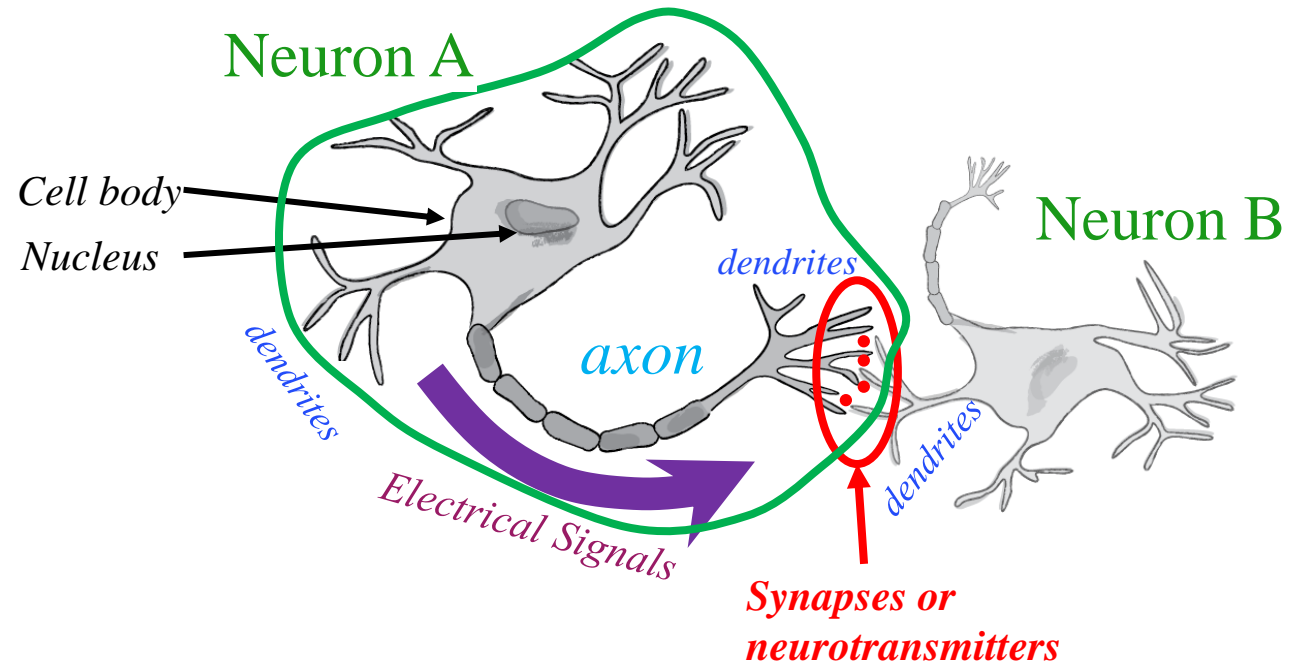❏ Neurons are connected together through *synapses* (~10K)



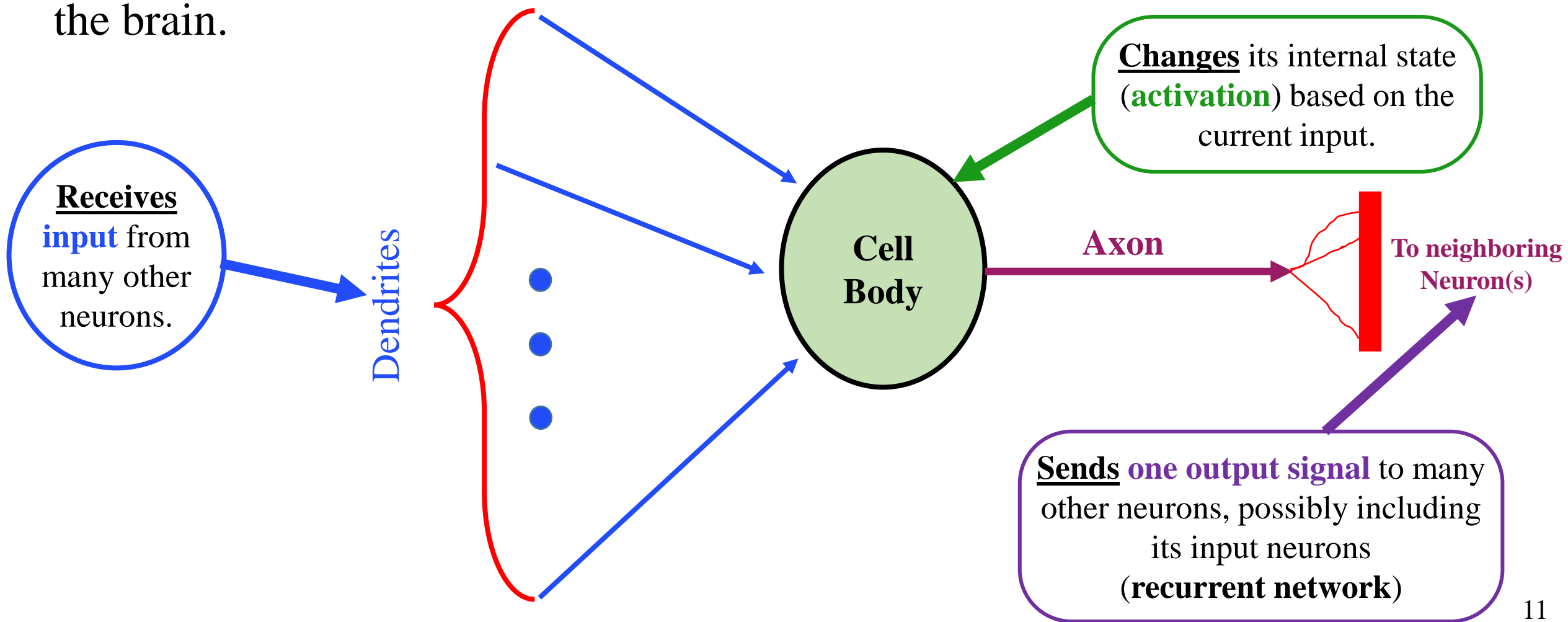**Image source:** https://beautifulnow.is/discover/wellness/new-brain-flows-are-beautiful-now



**Image source:** https://www.getbodysmart.com/nervous-system/neuron-synapse-structure

9

❑ A neuron accept (**and combine**) inputs through *dendrites* from other neurons

❑ If a given neuron *combined* input above a **threshold**, the neuron discharges a spike (**electrical pulse**) that travels from the body, down the **axon**, to the next **neuron(s)**

❑ The strength of the signal that reaches the next neuron depends on factors such as the amount of neurotransmitter (*synapses*) available

Neuron A

Cell body

Nucleus

dendrites

axon

Electrical Signals

dendrites

Neuron B

dendrites

Synapses or neurotransmitters

https://natureofcode.com/book/chapter-10-neural-networks/
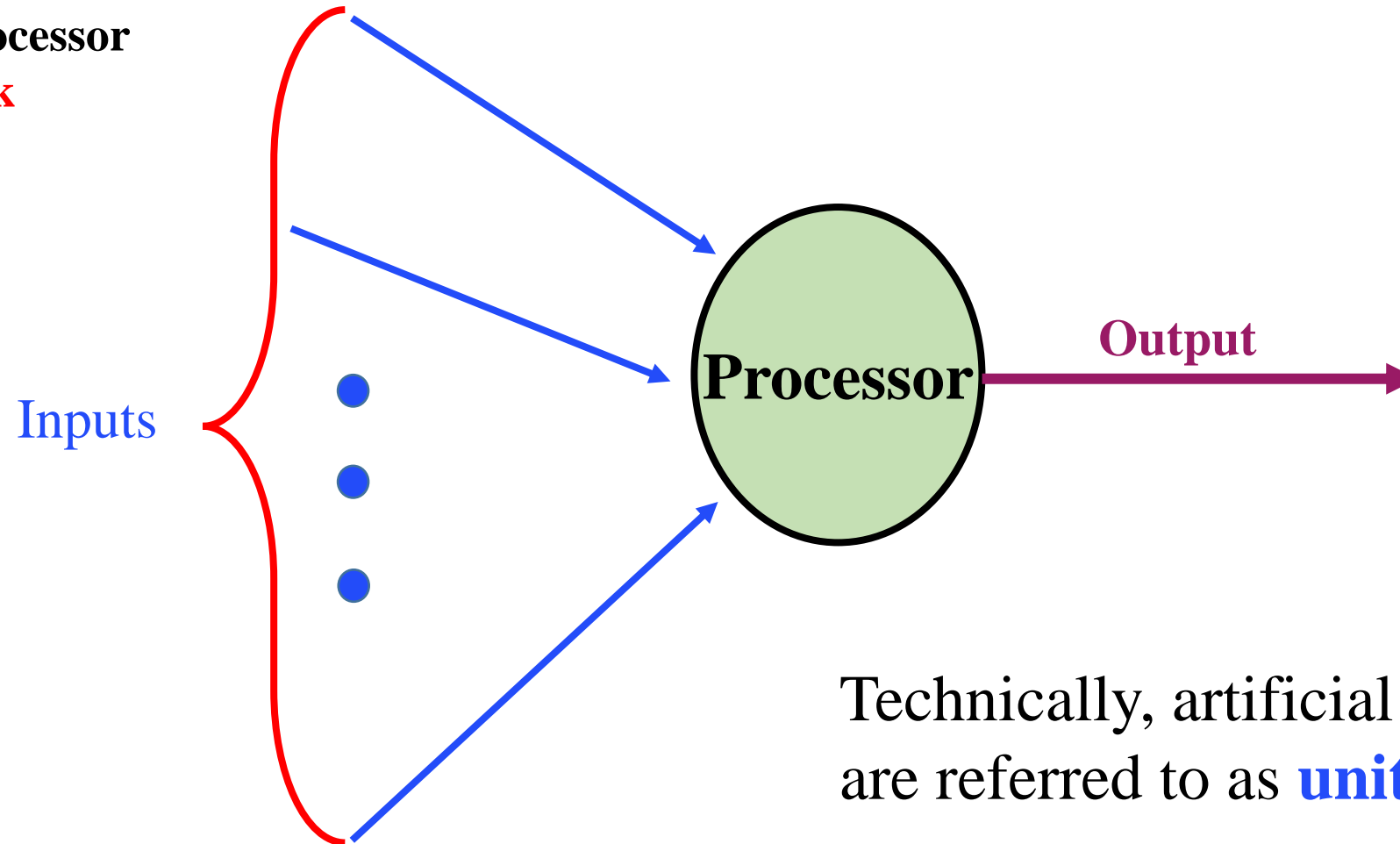
# Modeling of a Biological Neuron

❑ A mathematical model of the neuron (called the perceptron) has been introduced in an effort to mimic our understanding of the functioning of the brain.



**Receives** **input** from many other neurons.

Dendrites

**Changes** its internal state (**activation**) based on the current input.

**Cell Body**

**Axon**

To neighboring Neuron(s)

**Sends** **one output signal** to many other neurons, possibly including its input neurons (**recurrent network**)

# Artificial Neuron

❑ An artificial neuron is an imitation of a human neuron

➢ **Dendrites: Input**
➢ **Cell body: Processor**
➢ **Synaptic: Link**
➢ **Axon: Output**

Inputs

**Processor**

**Output**

Technically, artificial neurons are referred to as **units** or **nodes**.

12

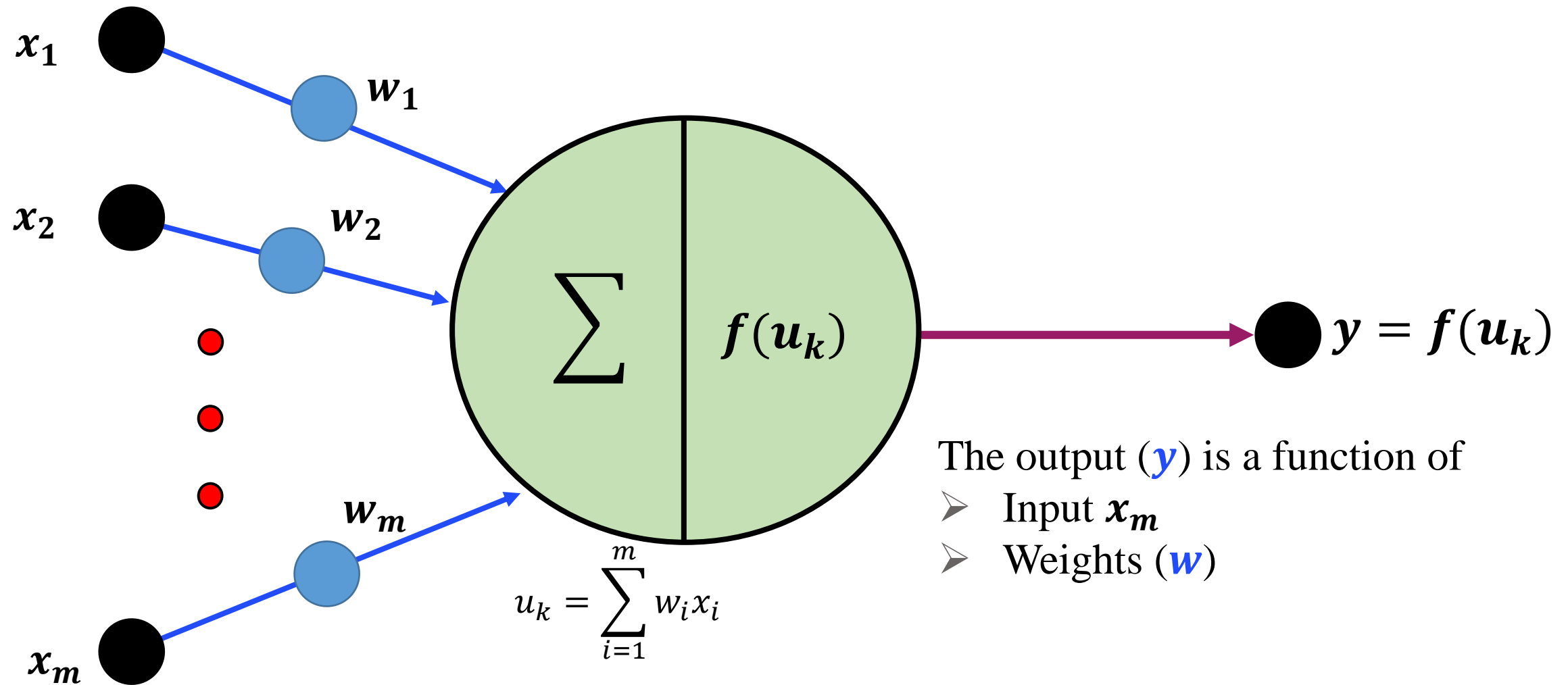Multiple inputs ($x$) each of which has a different strength, i.e., a **weight $w$**



$x_1$
$w_1$

$x_2$
$w_2$

Inputs

$w_m$
$x_m$

**Activation** the combined input must be above certain threshold

Sum Unit

Activation

**Output = $y$**

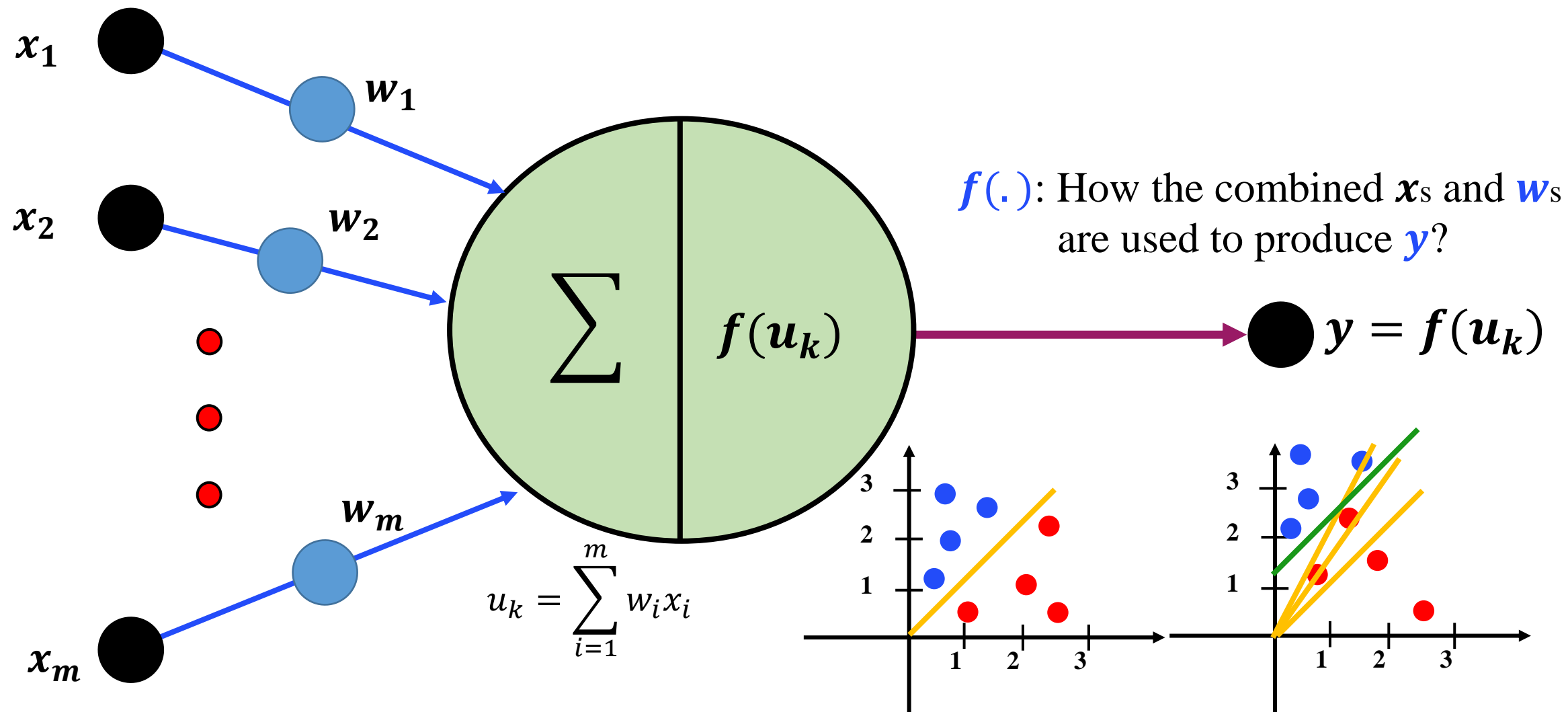**The operations done by a Neuron are:**
1) *Multiply* inputs by the weights,
2) *Add* them up
3) *Check* the sum against the activation and get y

Technically, artificial neurons are referred to as **units** or **nodes**.

13

# Artificial Neuron (cont'd)



$x_1$

$w_1$

$x_2$

$w_2$

$w_m$

$x_m$

$\sum$

$f(u_k)$

$u_k = \sum_{i=1}^{m} w_i x_i$
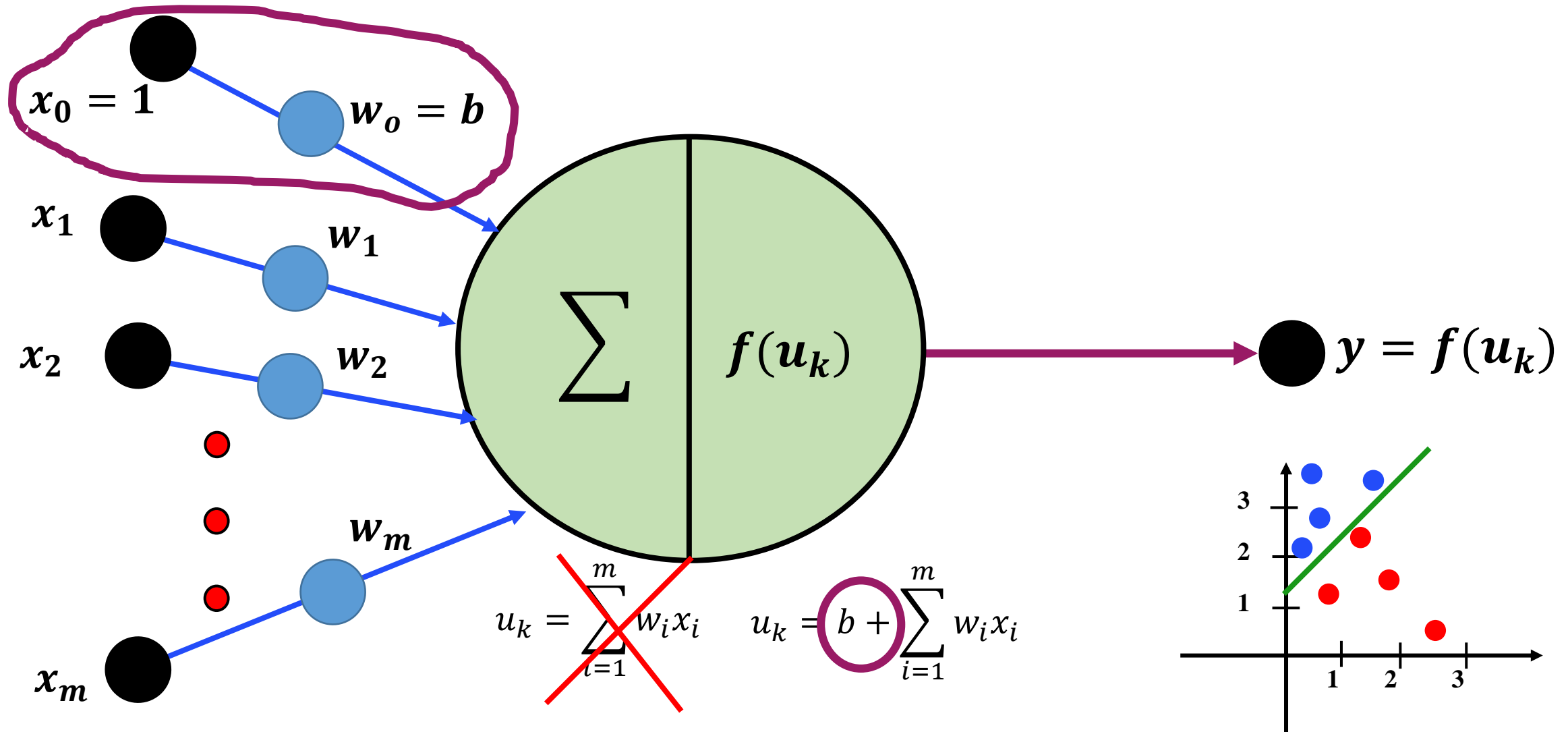
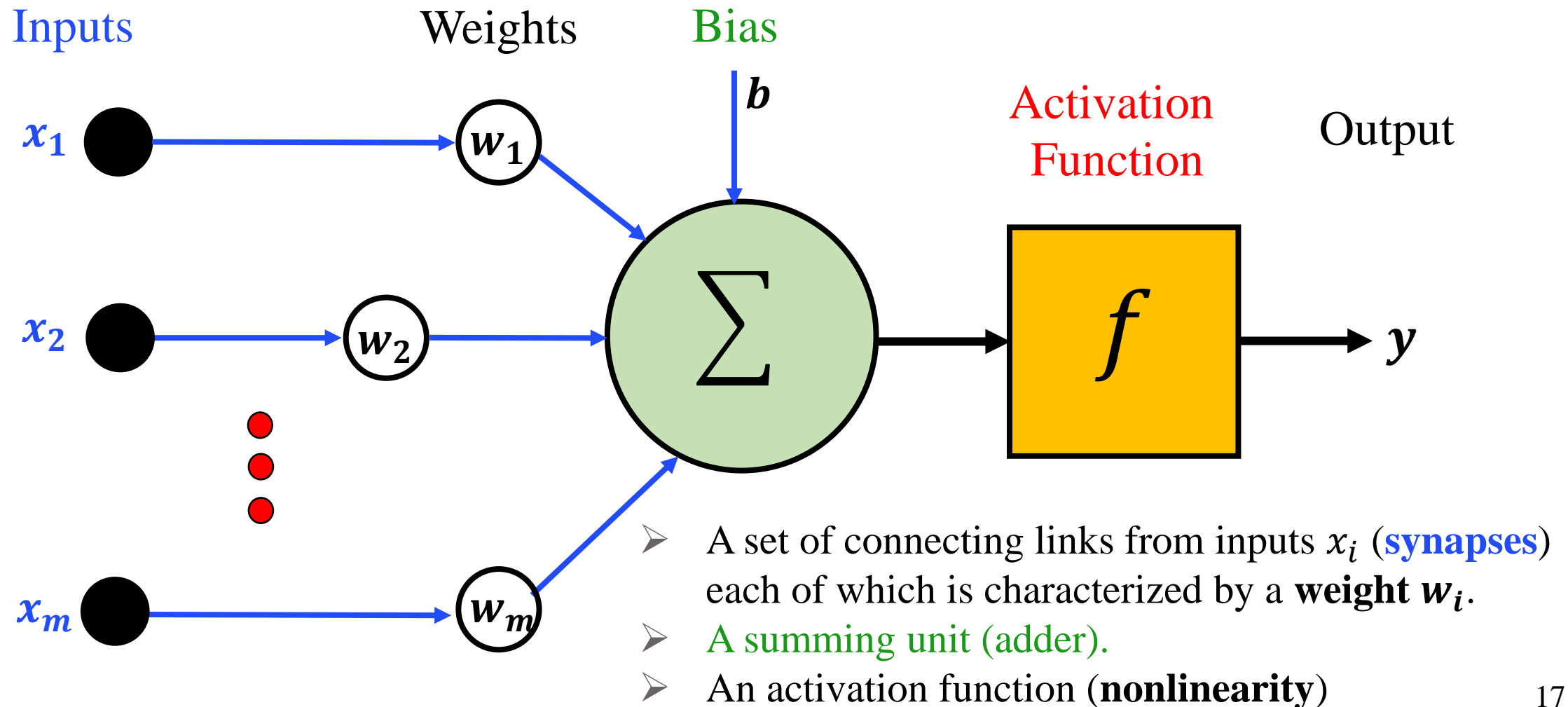$f(.)$: How the combined $x$s and $w$s are used to produce $y$?

$y = f(u_k)$

A **bias value** (**b**) is important to **full control** of the **activation function** (i.e., the output) for successful learning. This is a sort of **regularization**

$$x_0 = 1 \qquad w_o = b$$
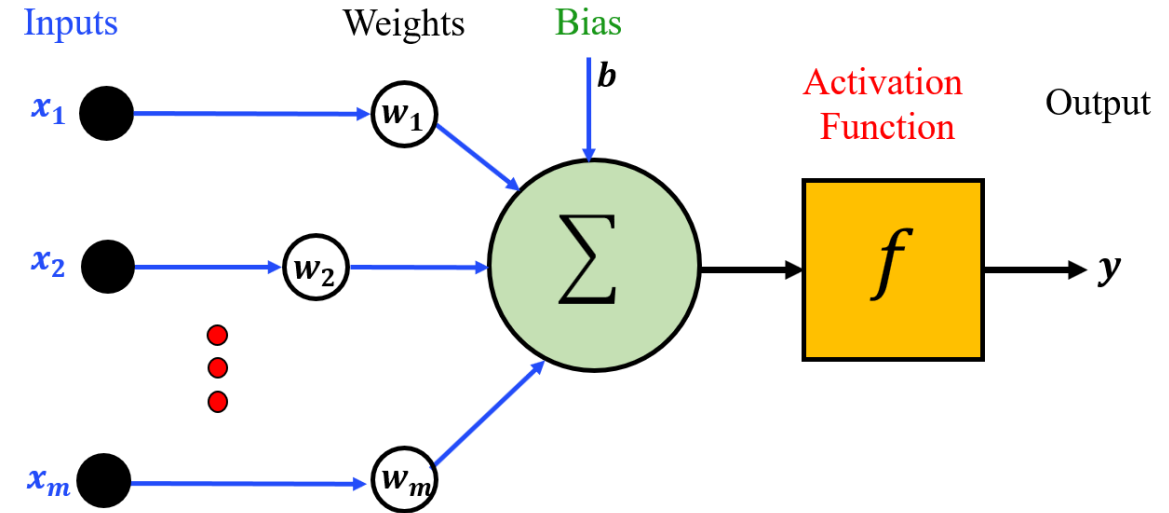
$$x_1 \qquad w_1$$

$$x_2 \qquad w_2$$

$$w_m$$

$$x_m$$

$$\sum \qquad f(u_k)$$

$$y = f(u_k)$$

$$u_k = \cancel{\sum_{i=1}^{m} w_i x_i} \qquad u_k = \left( b + \right)\sum_{i=1}^{m} w_i x_i$$

# Artificial Neuron Network (ANN)

Basic Elements of any ANN:



> A set of connecting links from inputs $x_i$ (**synapses**) each of which is characterized by a **weight $w_i$**.

> A summing unit (adder).

> An activation function (**nonlinearity**)

# ANN (cont'd)

- ❑ If the sum exceeds a certain threshold, the ANN (or the *perceptron*) fires an output value that is transmitted to the next unit(s)
- ❑ ANN uses nonlinear transfer function

**Why do we need nonlinearity?**

Inputs   Weights   Bias   Activation Function   Output

$$y = f\left(b + \sum_{i=1}^{m} w_i x_i\right) \implies y = f(b + \mathbf{W^T X})$$

y is **linear** and **unbounded**
- ➤ NOT realistic
- ➤ Can NOT be generalized
- ➤ LESS power to solve *complex nonlinear* problems
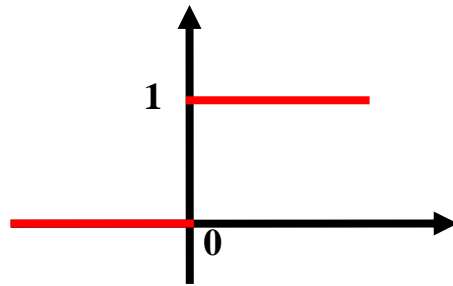
# ANN Transfer Functions

**Linear**

$$y_k = u_k$$

**Saturating linear**

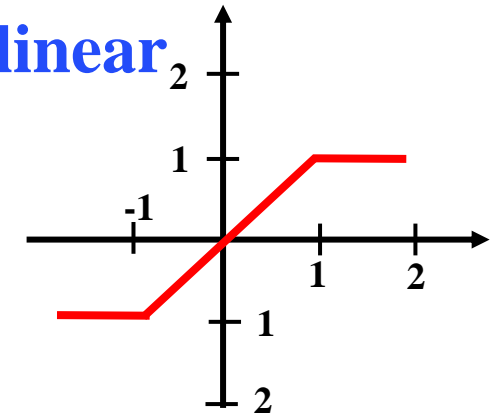$$y_k = \begin{cases} 1 & if\ u_k > 1 \\ u_k & if\ 0 \leq u_k \leq 1 \\ 0 & if\ u_k < 0 \end{cases}$$

**Hard Limit**

$$y_k = \begin{cases} 1 & if\ u_k \geq 0 \\ 0 & if\ u_k < 0 \end{cases}$$

**Symmetric Saturating linear**

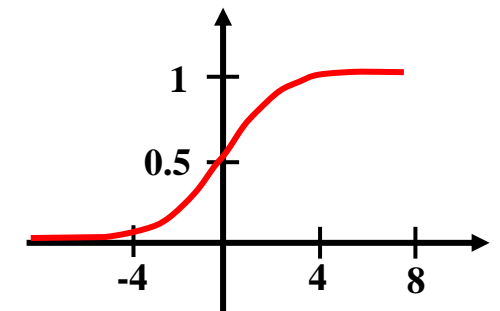$$y_k = \begin{cases} 1 & if\ u_k > 1 \\ u_k & if\ 0 \leq u_k \leq 1 \\ -1 & if\ u_k < 0 \end{cases}$$

**Symmetric Hard Limit**

$$y_k = \begin{cases} 1 & if\ u_k \geq 0 \\ -1 & if\ u_k < 0 \end{cases}$$
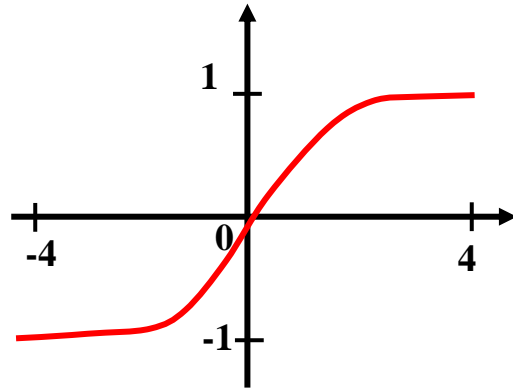
**Log Sigmoid**

$$y_k = \frac{1}{1 + e^{-u_k}}$$

19

# Artificial Neuron: Transfer Function
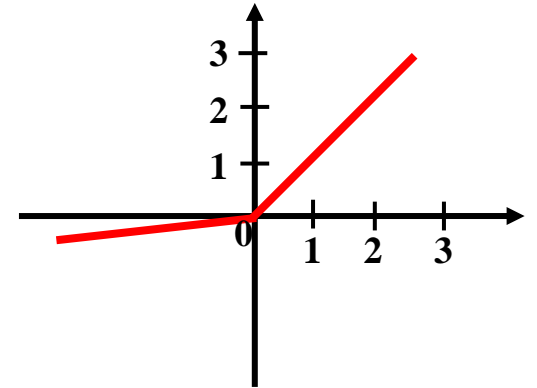
## Hyperbolic Tangent Sigmoid

$$y_k = \frac{e^{u_k} - e^{-u_k}}{e^{u_k} - e^{-u_k}}$$
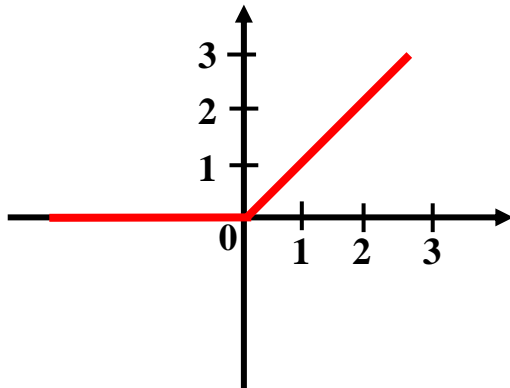


## Leaky ReLU
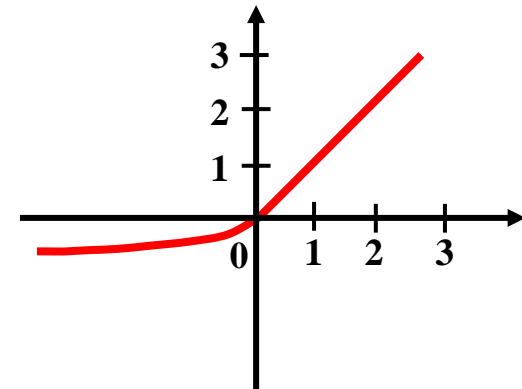
$$y_k = \max(\epsilon u_k, u_k)$$
$$\epsilon \ll 1$$


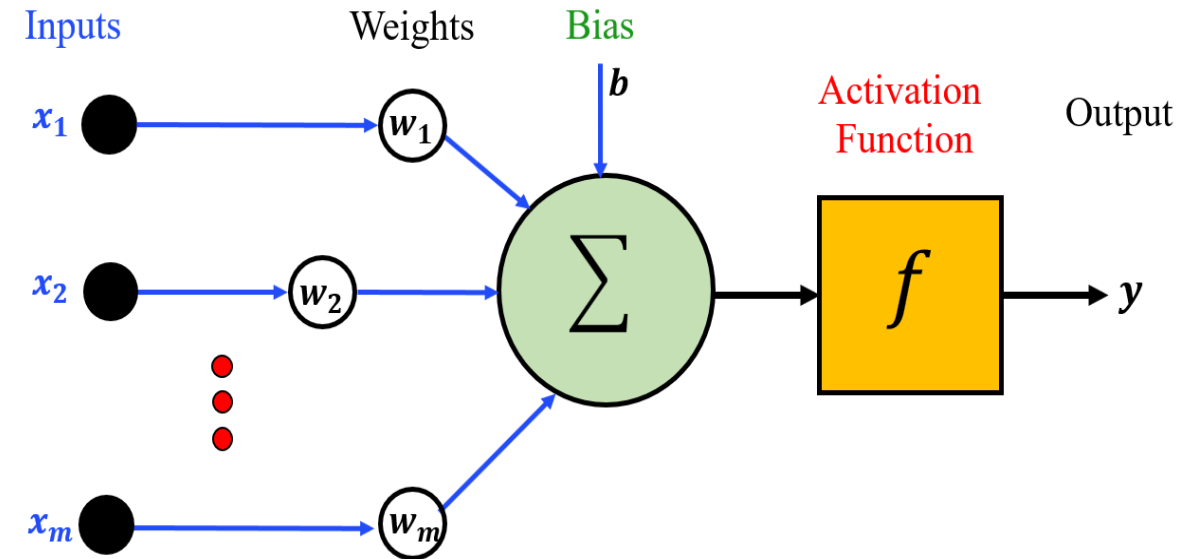
## Rectified Linear Unit (ReLU)

$$y_k = \max(0, u_k)$$



## Exponential Linear Unit (ELU)

$$y_k = \begin{cases} S_k & if\ u_k \geq 0 \\ \alpha(e^{S_k} - 1) & if\ u_k < 0 \end{cases}$$

# Artificial Neural Network (ANN)

❑ An artificial neural network (ANN) is a **massively parallel distributed processor** made up of **simple** processing units (neurons).

❑ ANN is capable of resolving paradigms that linear computing cannot resolve.

❑ ANNs are **adaptive systems**, i.e., parameters **can be changed** through a *learning* process (**training**) to suit the underlying problem.



❑ ANNs can be used in a *wide* variety of classification tasks, e.g., character recognition, speech recognition, fraud detection, medical diagnosis.

❑ *"neural networks are the second-best way of doing just about anything"* ***John Denker (AT&T Bell laboratories)***

21

# Learning Process

❑ **learning** is the process by which the *parameters* of an ANN, i.e., *w*, are **adapted** through a process of stimulation by the environment by which the network is embedded.

## Learning ≡ Training

➢ **Selection** of the network topology

➢ **Adapt** weights values.
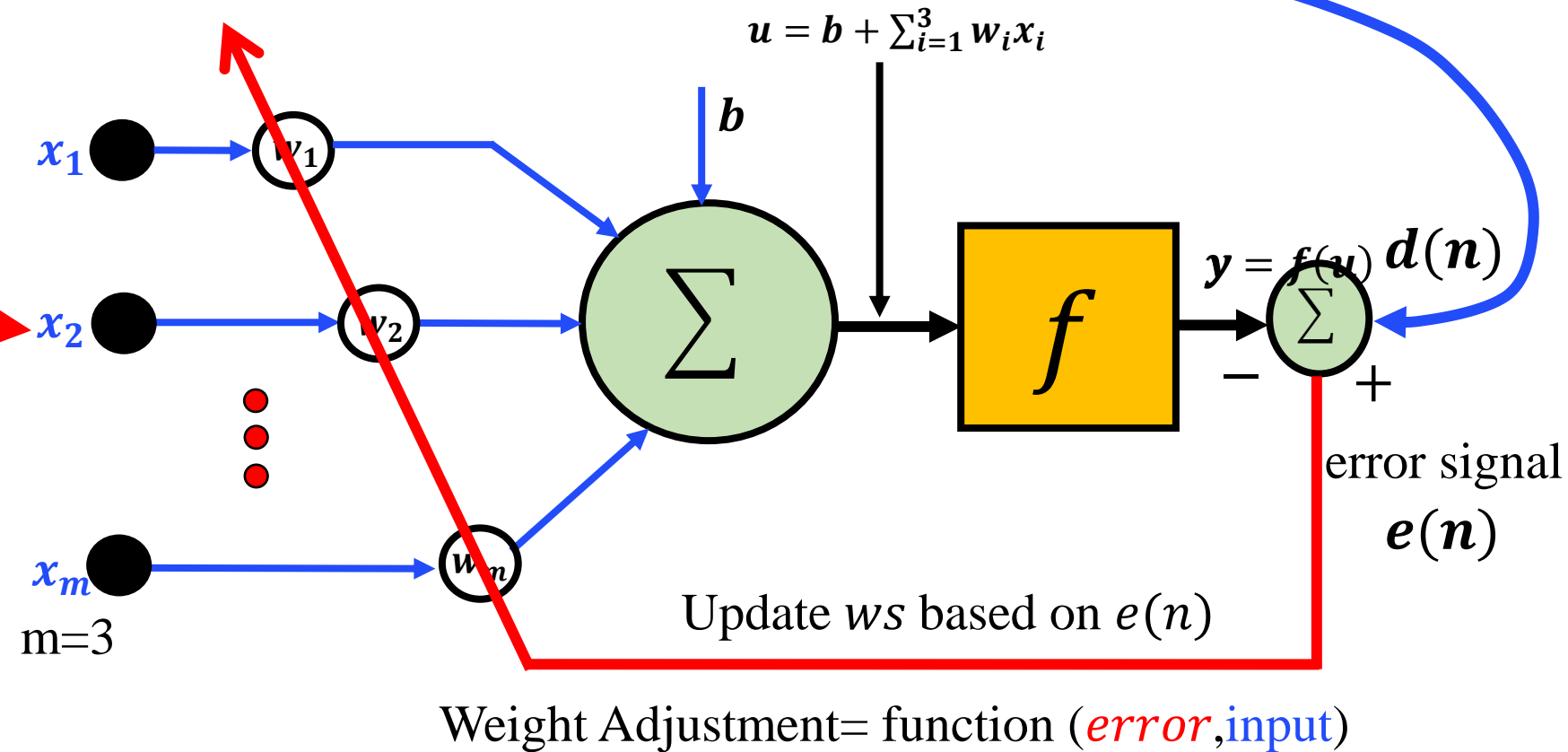
➢ **Learn** by trial-and-error (experience!)

❑ **Every data sample** for an ANN *training* consists of a vector **X(n)** and the corresponding (*desired* or *target*) output **d**

❑ A **batch** is a group of input samples with their *desired* outputs

| n | $x_1$ | $x_2$ | $x_3$ | Output | *d* Target |
|---|---|---|---|---|---|
| 1 | 10.33 | 56 | 0.56 | 0.8 | |
| 2 | 8.97 | 48 | 0.61 | 0.1 | |
| 3 | 11.01 | 49 | 0.49 | 0.3 | |
| 4 | 9.32 | 53 | 0.89 | 0.7 | |
| 5 | 10.51 | 50 | 0.71 | 0.4 | |
| 6 | 12.10 | 59 | 0.90 | 0.8 | |
| ⋮ | | | | | |
| 1996 | 7.99 | 61 | 0.59 | 0.9 | |
| 1997 | 11.36 | 52 | 0.63 | 0.5 | |
| 1998 | 12.09 | 48 | 0.78 | 0.2 | |
| 1999 | 10.81 | 55 | 0.87 | 0.7 | |
| 2000 | 13.00 | 53 | 0.91 | 0.6 | |

X(n) *Sample Features*

*Sample number*

*batch*

22

| n | $x_1$ | $x_2$ | $x_3$ | Output |
|---|---|---|---|---|
| 1 | 10.33 | 56 | 0.56 | 0.8 |
| 2 | 8.97 | 48 | 0.61 | 0.1 |
| 3 | 11.01 | 49 | 0.49 | 0.3 |
| 4 | 9.32 | 53 | 0.89 | 0.7 |
| 5 | 10.51 | 50 | 0.71 | 0.4 |
| 6 | 12.10 | 59 | 0.90 | 0.8 |
| ⋮ | | | | |
| 1996 | 7.99 | 61 | 0.59 | 0.9 |
| 1997 | 11.36 | 52 | 0.63 | 0.5 |
| 1998 | 12.09 | 48 | 0.78 | 0.2 |
| 1999 | 10.81 | 55 | 0.87 | 0.7 |
| 2000 | 13.00 | 53 | 0.91 | 0.6 |

$$u = b + \sum_{i=1}^{3} w_i x_i$$

$y = f(u)$  $d(n)$

error signal $e(n)$

$x_1$, $x_2$, $x_m$  m=3

Update $ws$ based on $e(n)$

Weight Adjustment = function ($error$, input)

23

| n | $x_1$ | $x_2$ | $x_3$ | Output |
|---|---|---|---|---|
| 1 | 10.33 | 56 | 0.56 | 0.8 |
| 2 | 8.97 | 48 | 0.61 | 0.1 |
| 3 | 11.01 | 49 | 0.49 | 0.3 |
| 4 | 9.32 | 53 | 0.89 | 0.7 |
| 5 | 10.51 | 50 | 0.71 | 0.4 |
| 6 | 12.10 | 59 | 0.90 | 0.8 |
| ⋮ | | | | |
| 1996 | 7.99 | 61 | 0.59 | 0.9 |
| 1997 | 11.36 | 52 | 0.63 | .5 |
| 1998 | 12.09 | 48 | 0.78 | .2 |
| 1999 | 10.81 | 55 | 0.87 | 0.7 |
| 2000 | 13.00 | 53 | 0.91 | 0.6 |

new input sample(s) → output → update weights

Update $ws$ based on $e(n)$

Weight Adjustment= function (*error*,input)

**General rule for neuron learning**

$$w_{new} = w_{old} + \eta * e * x$$

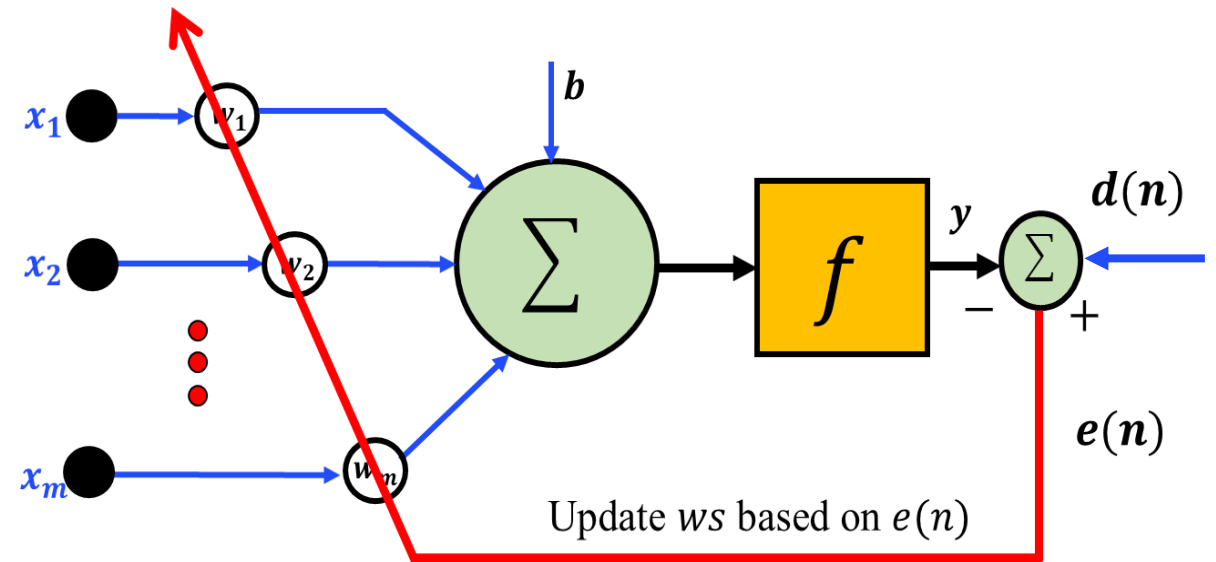$\eta$ is the learning constant or the ***learning rate***

24

❑ **Learning** is a *recursive* operation through which network **parameters** (*weights*) are *updated* in a way to reduce the **difference** (*error*) between network output and the *desired* (**target**) output

**Set** initial values of the weights (e.g., randomly)

*Do*

    **Compute** the output function of a given input ($X(n)$)

    **Evaluate** the output by comparing $y(n)$ with $d(n)$.

    **Adjust** the *weights*.

    **Loop** until a criterion is met.

*end*



Update *ws* based on $e(n)$

## Criterion

➢ Certain number of iterations

➢ Error threshold

# Learning Process: Cost Function

- Our **objective** is to **reduce** the difference between the *actual* and *target* outputs (i.e., the error)

- This can be achieved by **minimizing** a **function** of the error (**error energy**)
  - This called the *cost function*.
  - Example is the **mean squared error**



Update *ws* based on $e(n)$

$$e(n) = d(n) - y(n)$$

$$E(n) = \frac{1}{2} e^2(n) = \frac{1}{2} \left( d(n) - y(n) \right)^2$$

- This learning is called *error-correction learning* or *delta* rule or *Widrow-Hoff* rule

$$\Delta w_{kj}(n) = \eta . e_k(n) . x_j(n)$$

$n$ is the current sample
$k$ index for the current neuron
$j: 1 \rightarrow m$

$$w_{kj}(n+1) = w_{kj}(n) + \Delta w_{kj}(n)$$

- The adjustment of a weight vector of **_n input neuron connection_** is proportional to the *product* of the **error signal** and the **input value** of the connection in question.

26

| n | $x_1$ | $x_2$ | $x_3$ | Output |
|---|---|---|---|---|
| 1 | 10.33 | 56 | 0.56 | 0.7 |
| 2 | 8.97 | 48 | 0.61 | 0.9 |
| 3 | 11.01 | 49 | 0.49 | 0.8 |
| 4 | 9.32 | 53 | 0.89 | 0.8 |
| 5 | 10.51 | 50 | 0.71 | 0.7 |
| 6 | 12.10 | 59 | 0.90 | 0.8 |
| ⋮ | | | | |
| 1996 | 7.99 | 61 | 0.59 | 0.9 |
| 1997 | 11.36 | 52 | 0.63 | 0.9 |
| 1998 | 12.09 | 48 | 0.78 | 0.8 |
| 1999 | 10.81 | 55 | 0.87 | 0.7 |
| 2000 | 13.00 | 53 | 0.91 | 0.6 |

Update $ws$ based on $e(n)$

The training cycle at which **All** the training samples have been used by the network is called the *epoch*

27

**Example**

| n | $x_1$ | $x_2$ | $x_3$ | d |
|---|-------|-------|-------|-----|
| 1 | 1 | 1 | 0.5 | 0.7 |
| 2 | -1 | 0.7 | -0.5 | 0.2 |
| 3 | 0.3 | 0.3 | -0.3 | 0.3 |

*Assume*
- initial weights are 0.5, -0.3, 0.8,
- b=0;
- $\eta$=0.1 and
- linear activation function

**Solution**

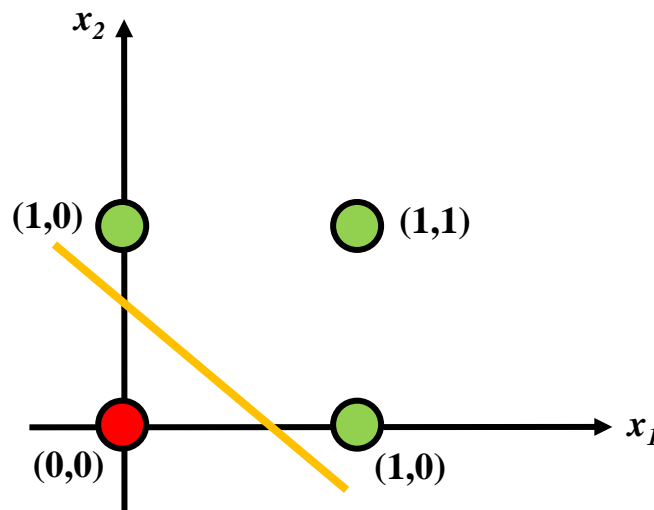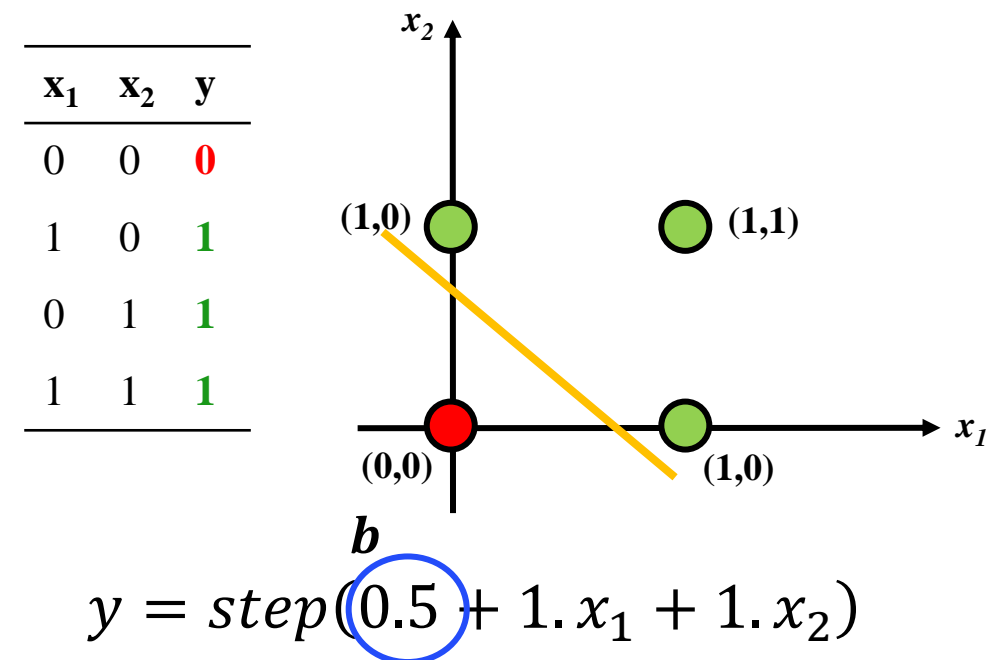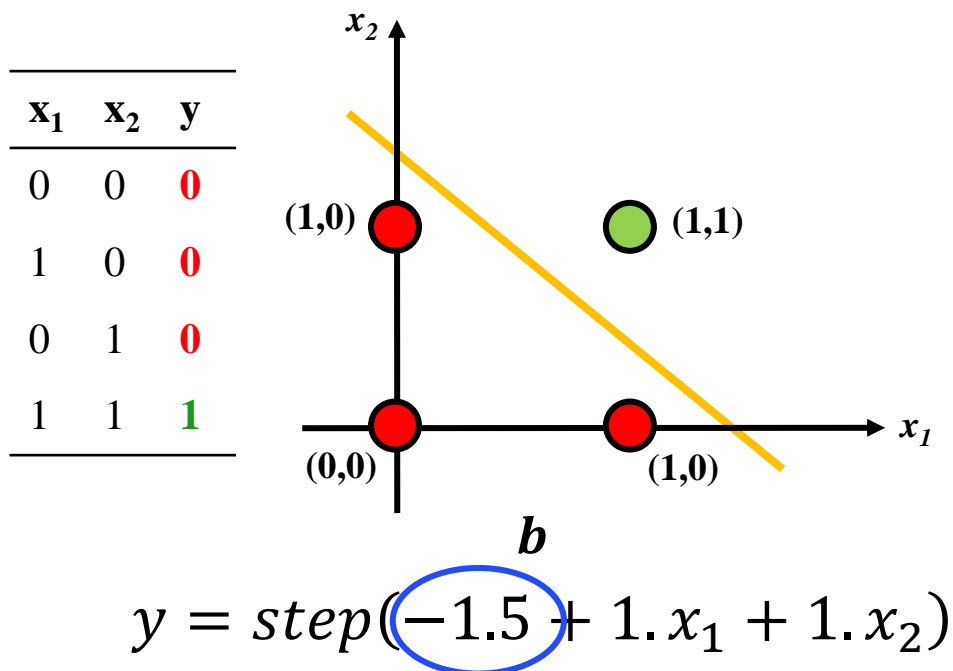| n | $x_1$ | $x_2$ | $x_3$ | d |
|---|-------|-------|-------|-----|
| 1 | 1 | 1 | 0.5 | 0.7 |
| 2 | -1 | 0.7 | -0.5 | 0.2 |
| 3 | 0.3 | 0.3 | -0.3 | 0.3 |

# ANN Examples

- One layer *feedforward* neural network called the *perceptron*

- Can solve linear function, e.g., AND, OR, NOT

| $x_1$ | $x_2$ | $y$ |
|-------|-------|-----|
| 0 | 0 | **0** |
| 1 | 0 | **0** |
| 0 | 1 | **0** |
| 1 | 1 | **1** |



$$y = f\left( b + \sum_{i=1}^{n} w_i x_i \right)$$

$$y = step(-1.5 + 1.x_1 + 1.x_2)$$

❑ One layer *feedforward* neural network called the *perceptron*
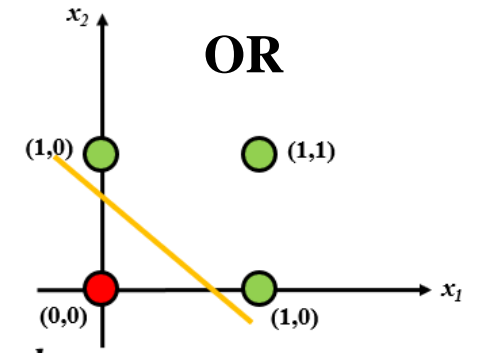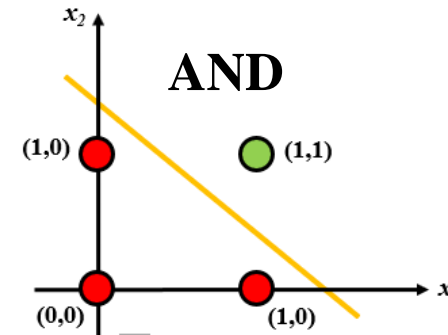
❑ Can solve linear function, e.g., AND, OR, NOT

| $x_1$ | $x_2$ | $y$ |
|-------|-------|-----|
| 0 | 0 | 0 |
| 1 | 0 | 1 |
| 0 | 1 | 1 |
| 1 | 1 | 1 |



$$y = f\left(b + \sum_{i=1}^{n} w_i x_i\right)$$

$$y = step(0.5 + 1.x_1 + 1.x_2)$$

31

# ANN Examples (cont'd)



Left table:

| $x_1$ | $x_2$ | y |
|---|---|---|
| 0 | 0 | 0 |
| 1 | 0 | 0 |
| 0 | 1 | 0 |
| 1 | 1 | 1 |

$$y = step(-1.5 + 1.x_1 + 1.x_2)$$

Right table:

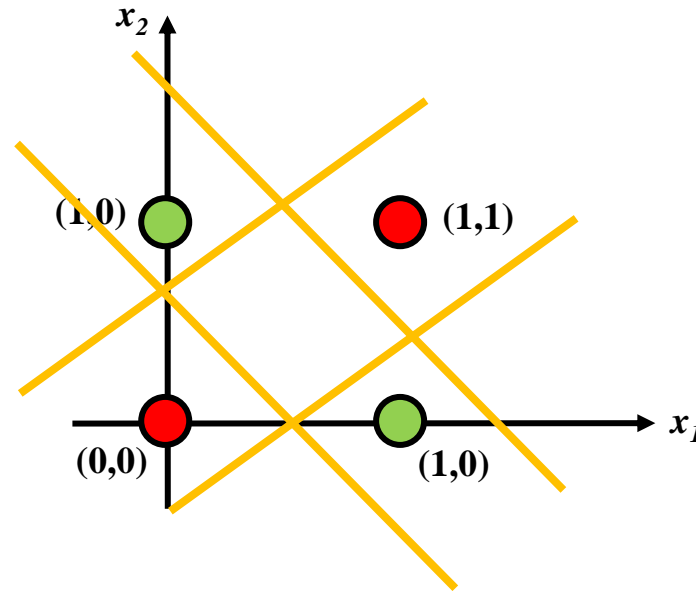| $x_1$ | $x_2$ | y |
|---|---|---|
| 0 | 0 | 0 |
| 1 | 0 | 1 |
| 0 | 1 | 1 |
| 1 | 1 | 1 |

$$y = step(0.5 + 1.x_1 + 1.x_2)$$

- ❑ Solving linearly, means the **decision boundary** is linear (straight line in 2D and a plane in 3D)
- ❑ The bias term (**b**) alters the **position**, but not the **orientation**, of the decision boundary
- ❑ The weights ($w_1$, $w_2$, ...$w_m$) determine the gradient
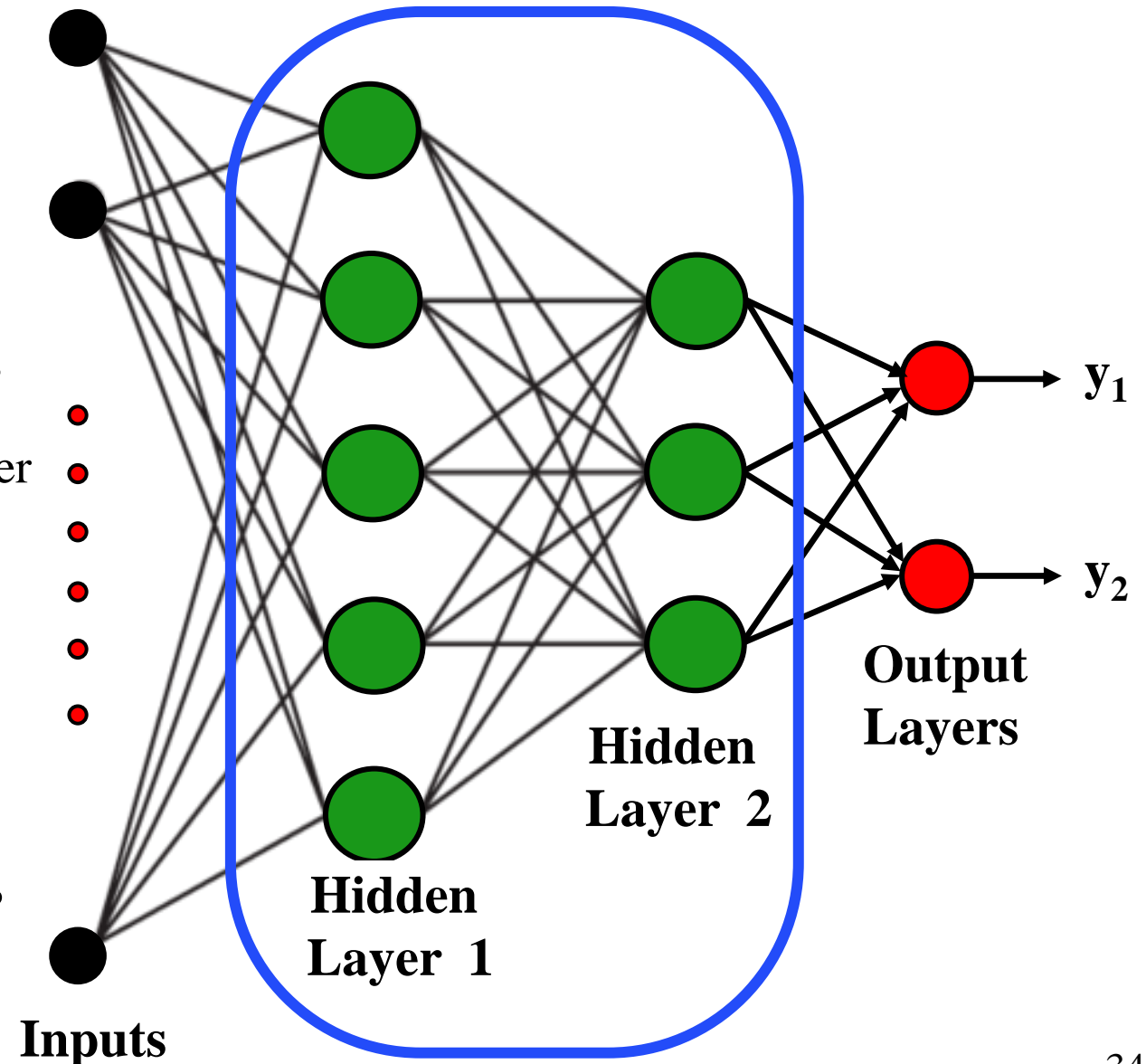
# ANN Examples: XOR function



| $x_1$ | $x_2$ | y |
|-------|-------|---|
| 0 | 0 | 0 |
| 1 | 0 | 1 |
| 0 | 1 | 1 |
| 1 | 1 | 0 |

- ❑ The XOR function is said to be **not linearly separable**
- ❑ If **one neuron** defines **one line** through input space, **what do we need to have two lines?**
- ❑ We need to have two neurons working in *parallel* (*next to each other rather than in different layers*).
- ❑ We would need a **multilayer neural network** to model (or to separate the two classes using) the XOR function.

# Multilayer Perceptron (MLP)

- More layers between the *input* the *output* layers

- Fully connected layers

- Multiple neurons at the output layers

  $y_j$ , $\mathbf{j} \in C$   C is set of all neurons at the output layer

- Error **backpropagation** is used for learning

  $$e(n) = d(n) - y(n)$$

- Weight adjustments are applied so as to minimize $e(n)$ in a statistical sense.



**Inputs**

**Hidden Layer 1**

**Hidden Layer 2**

**Output Layers**

$y_1$

$y_2$

# Gradient Descent

The **delta rule** is a gradient descent learning rule for updating the weights of an artificial neuron inputs in a single-layer NN

$$w_{kj}(n+1) = w_{kj}(n) + \Delta w_{kj}(n)$$



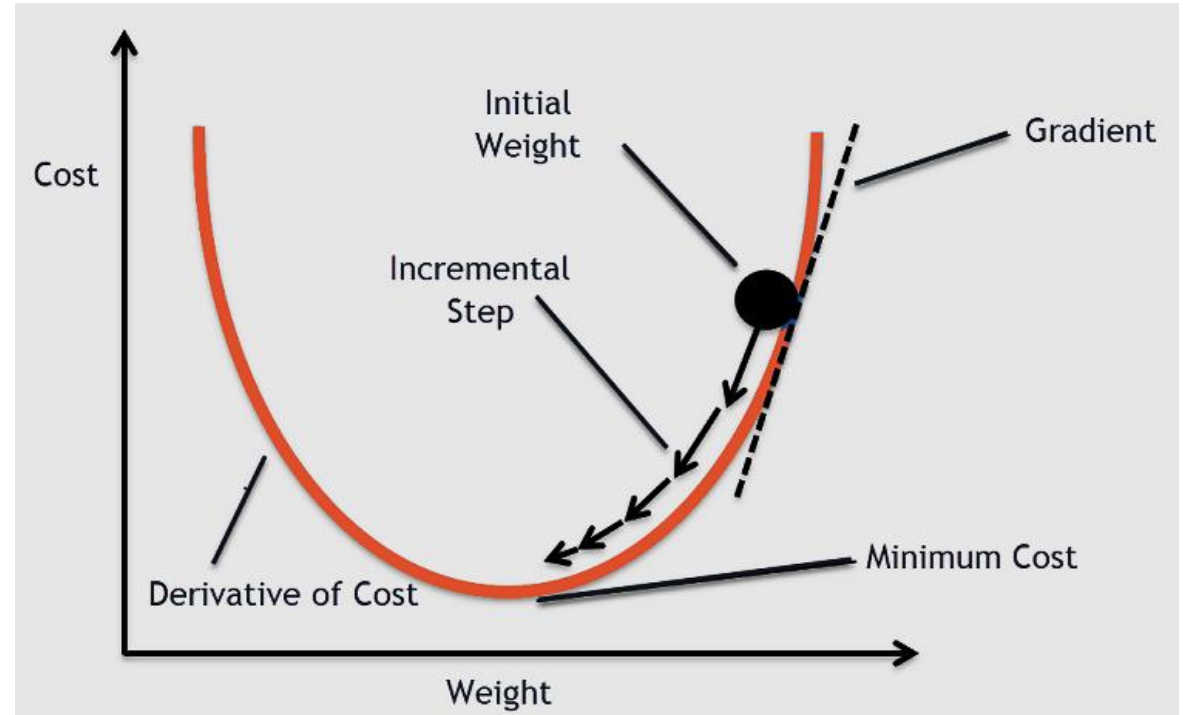**Image Source:** https://datascience-enthusiast.com/figures/cost.jpg



https://medium.com/@divakar_239/stochastic-vs-batch-gradient-descent-8820568eada1

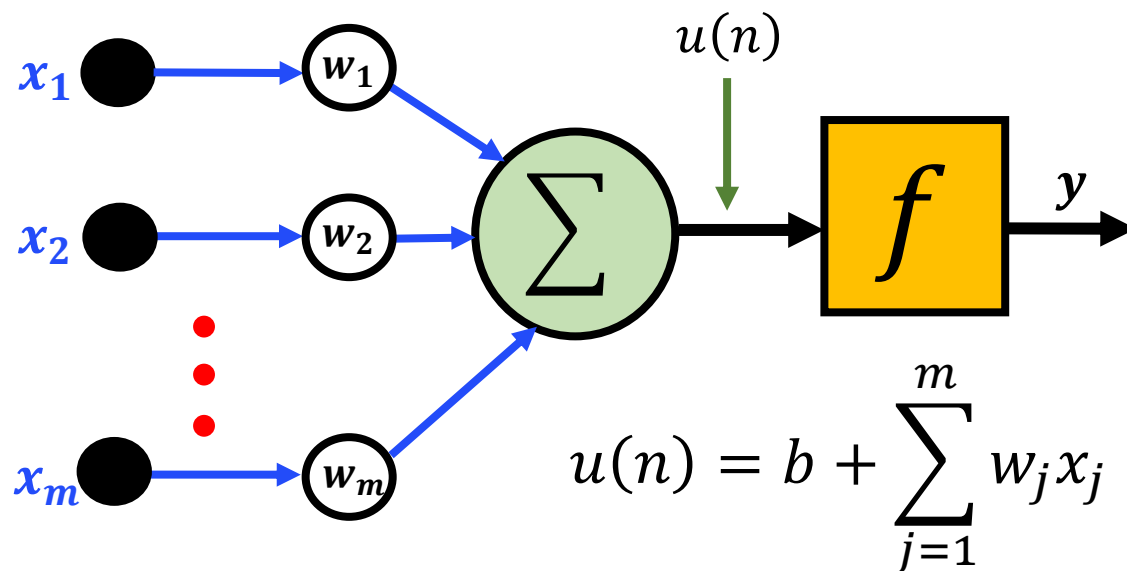The goal of gradient descent is to *iteratively* take steps towards **lower** regions (minima) of the loss function

# Gradient Descent (cont'd)

For *linear activation function*, the weight adjustment for a **neuron $k$** is given by
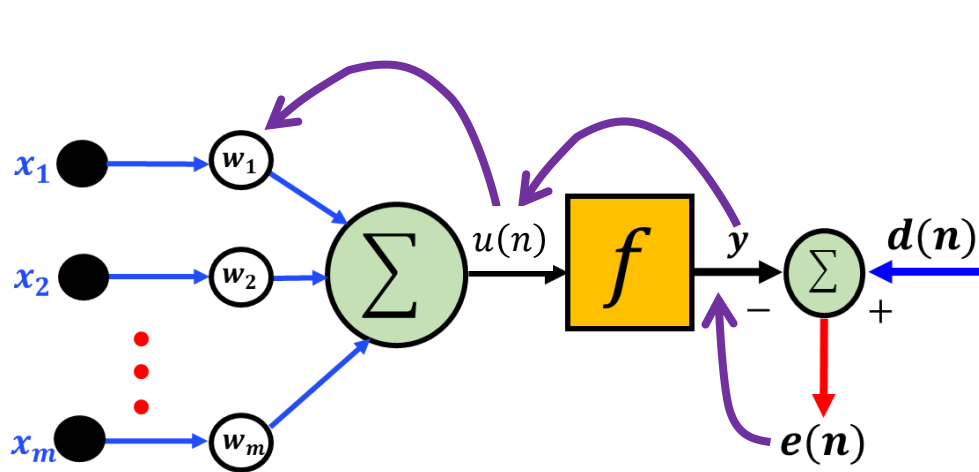
$$\Delta w_{kj}(n) = \eta * e_k(n) * x_j(n) \qquad\qquad j = 1,2,\dots m$$

For **any activation** function $f$:

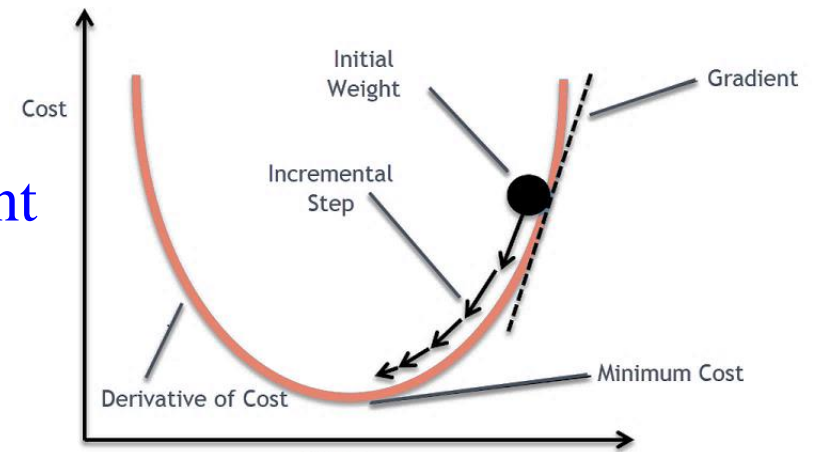$$\Delta w_{kj}(n) = \eta * e_k(n) * f'(u(n)) * x_j(n) *$$



$$u(n) = b + \sum_{j=1}^{m} w_j x_j$$

$$\Delta w_{kj} = \boxed{-} \eta * \frac{\partial E}{\partial w_j}$$

minimization    gradient

https://medium.com/@divakar_239/stochastic-vs-batch-gradient-descent-8820568eada1

By applying the chain rule

$$\frac{\partial E}{\partial w_j} = \left(\frac{\partial E}{\partial e}\right)\left(\frac{\partial e}{\partial y}\right)\left(\frac{\partial y}{\partial u}\right)\left(\frac{\partial u}{\partial w_j}\right)$$

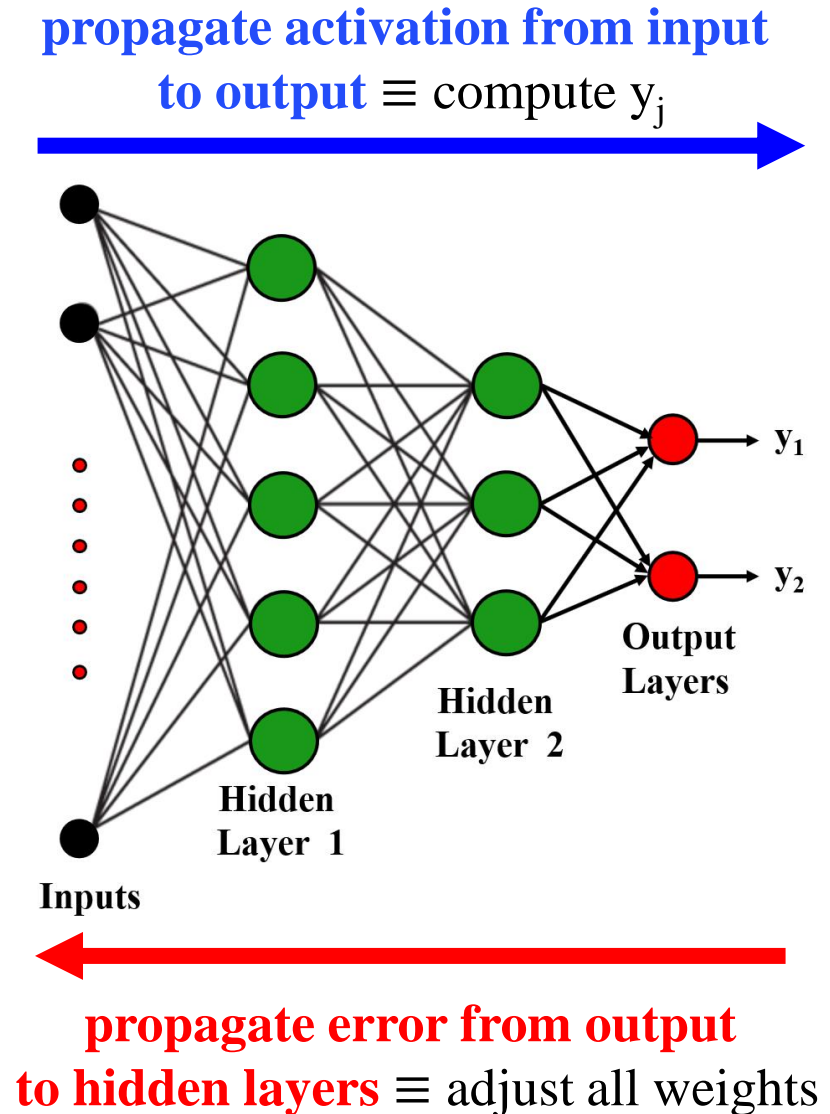$$\Delta w_{kj} = -\eta * (e)(-1)\big(f'(u(n))\big)(x_j)$$

$$\Delta w_{kj} = \eta * e * f'(u(n))x_j$$

$$E(n) = \frac{1}{2}e^2(n) \implies \frac{\partial E}{\partial e} = e$$

$$e(n) = d(n) - y(n) \implies \frac{\partial e}{\partial y} = -1$$

$$y(n) = f(u(n)) \implies \frac{\partial y}{\partial v} = f'(v(n))$$

$$u(n) = \sum_{j=1}^{m} w_j x_j \implies \frac{\partial u}{\partial w_j} = x_j$$

# **Backpropagation**

- ❑ Backpropagation is *supervised* algorithm that is a generalization for the least mean square (LMS) algorithm

- ❑ It is based on the ***gradient search*** technique to minimize the **cost function** ≡ squared error between the network output and the *target* output

- ❑ It is **recursive** application of the *chain rule* to compute the *gradients*

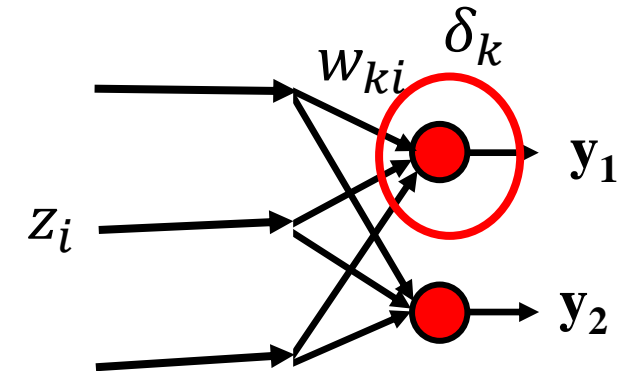Please see the following for all details about mathematical derivation: https://www.jeremyjordan.me/neural-networks-training/

**propagate activation from input to output** ≡ compute $y_j$



**propagate error from output to hidden layers** ≡ adjust all weights

38

❏ The weights of each **output neuron** can be determined directly using the *delta* learning rule.

$$\Delta w_{ki} = \eta * \boxed{e * f'(\cdot)} * z_i \qquad \delta_k = e * f'(\cdot)$$

**local gradient or error signal**

# Backpropagation (cont'd)

❑ The weights of each **output neuron** can be determined directly using the *delta* learning rule.
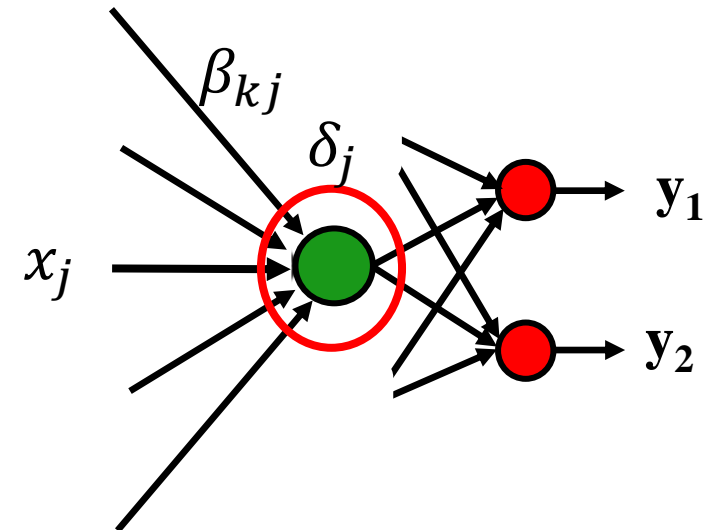
$$\Delta w_{ki} = \eta * \boxed{e * f'(\cdot)} * z_i \qquad \delta_k = e * f'(\cdot)$$

**local gradient or error signal**



❑ If the neuron is a *hidden* node

$$\delta_j = f'(\cdot) * \sum_{k=1}^{K} \delta_k * w_k$$

K is the set of all <u>nodes</u> on a ***next layer*** connected to the **current neuron**

[local gradient] x [upstream gradient]

Please see the following for all details about mathematical derivation:
https://www.jeremyjordan.me/neural-networks-training/

# Backpropagation Example

❑ Assume **one** input layer, **one** hidden layer, and **one** output neuron

$x_j$ :is the $j^{th}$ **input**

$z_i$ :is the output of the $i^{th}$ **hidden neuron**

$y_k$ :is the output of the $k^{th}$ **output neuron**

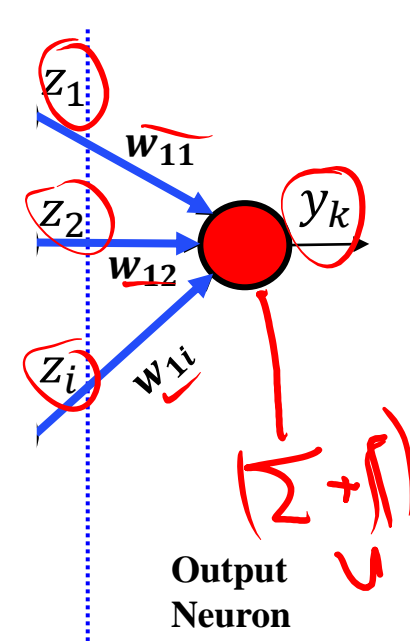$\beta_{ij}$ :is the weight from input node $x_j$ to hidden node $z_i$

$w_{ki}$ :is the weight from **hidden node** $z_i$ to **output neuron** $y_k$



**Output Neuron**

❑ The weights of the output neuron can be adjusted using the *delta learning* rule and the error signal:

$$\delta_{y_k} = e_k * f'(u_k) = (d_k - y_k)\, f'(u_k) \qquad u_k = \sum_{i=1}^{I} w_{ki} z_i$$

❑ Update the weights as follows:

$$w_{ki}(n+1) = w_{ki}(n) + \eta * \delta_{yk} * z_i$$

# Backpropagation (cont'd)

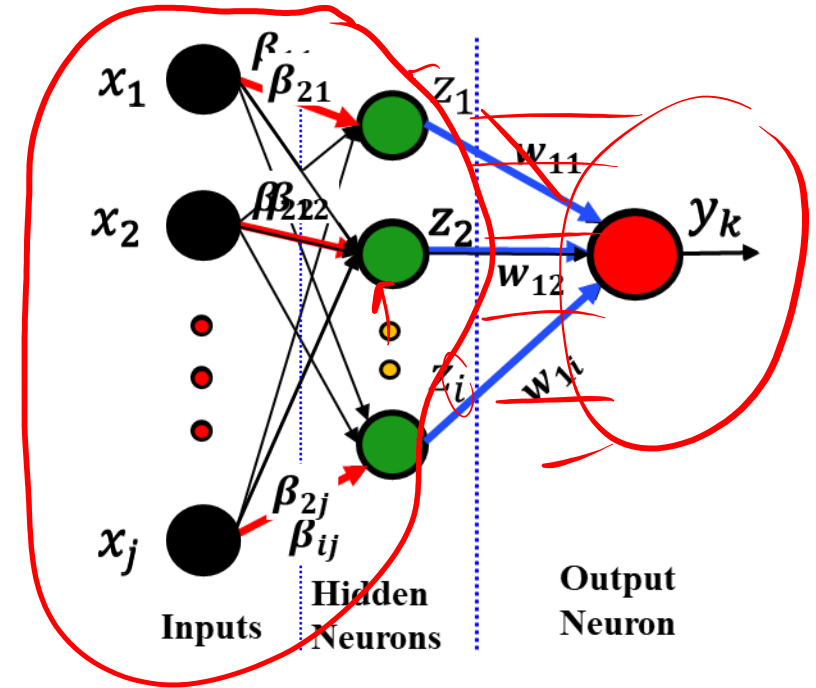- The weights of a $i^{th}$ hidden neuron can be adjusted using its error signal:

$$\delta_{z_i} = f'(u_i) * \sum_{k=1}^{K} \delta_{y_k} * w_{ki} \qquad u_i = \sum_{j=1}^{J} \beta_{ij} x_j$$



- Using the error signals, the weights of the $i^{th}$ hidden neuron can be updated

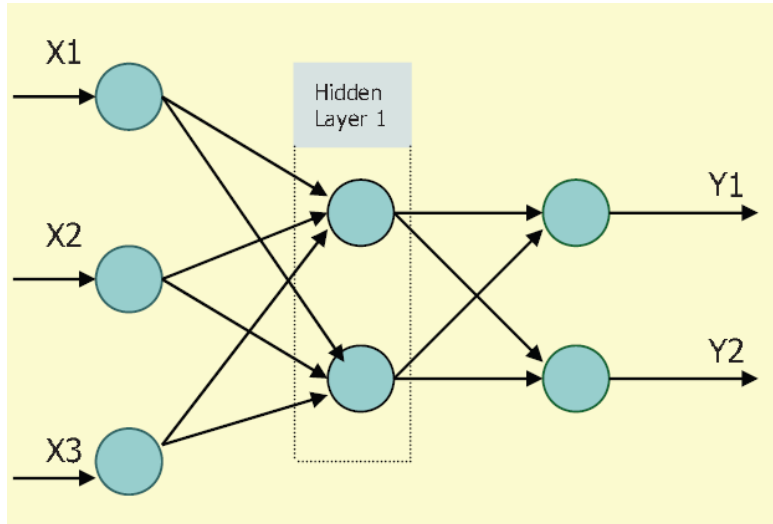$$\beta_{ij}(n+1) = \beta_{ij}(n) + \eta * \delta_{zi} * x_j$$

- For a sigmoid activation with zero bias
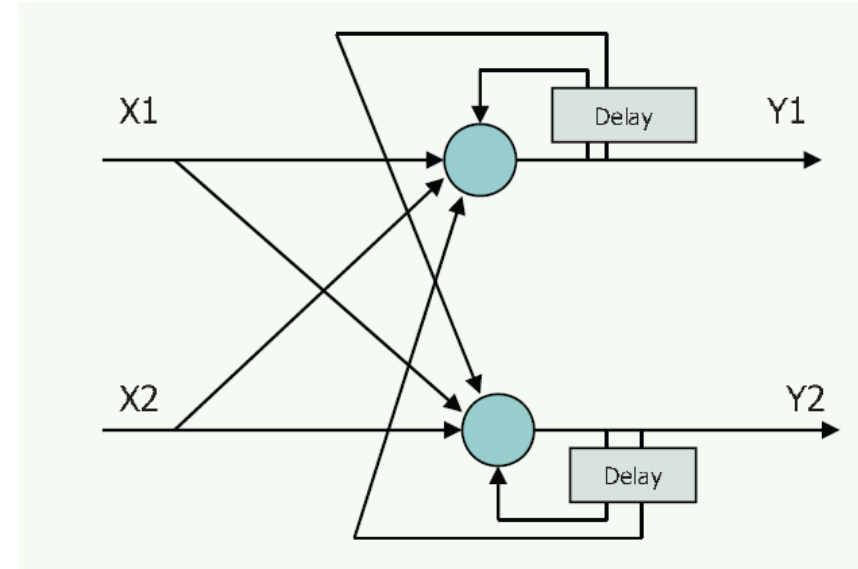
$$f'(u_k) = y_k(1 - y_k)$$

# Types of Neural Networks

**Feedforward** neural network



Signals to travel one way only (input to output)

Learning **with** a teacher

Supervised Learning

**Recurrent** neural network (RNN)



Output from previous step is fed as input in the current step
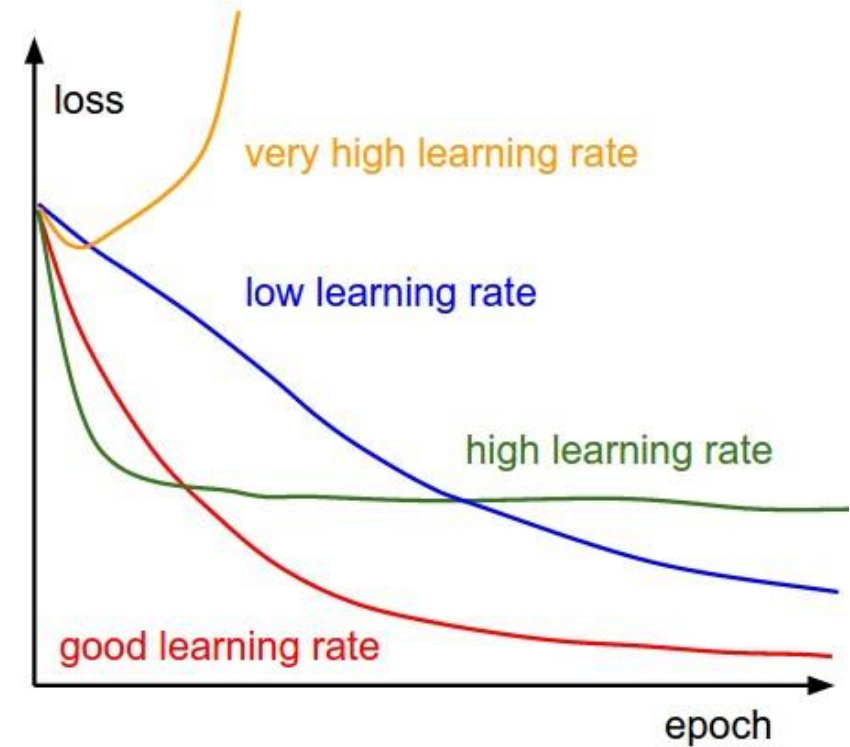
Learning **without** a teacher

Unsupervised Learning

Self-organizing maps (SOM)

# ANN Design and Issues

❑ Number of neurons, and hidden layers

❑ Initial weights (small random values $\in[-1,1]$)

❑ Choice of the transfer function

❑ Learning rate

❑ Weights adjusting
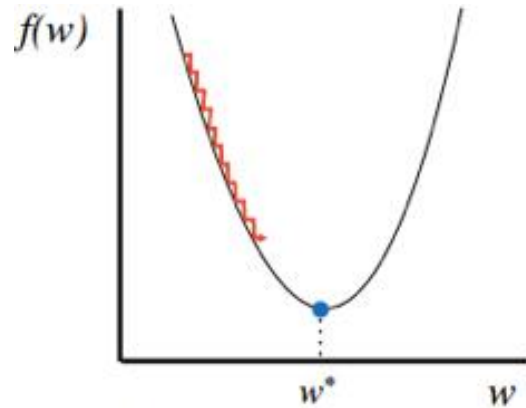
❑ Data representation, pre-processing, and splitting

# Learning Rate

❑ The learning rate, $\eta$, is a configurable (**hyper**)**parameter** used in ANNs training

❑ $\eta$ controls how *quickly* the model is adapted to the problem

❑ Practical value $0 < \eta < 1$.

> **Smaller** $\eta$ → smaller changes to $w$ → more training epochs

- Can cause the local minima stuck.

> Larger $\eta$ → larger changes to $w$ → fewer training epochs.

- May results in divergence.



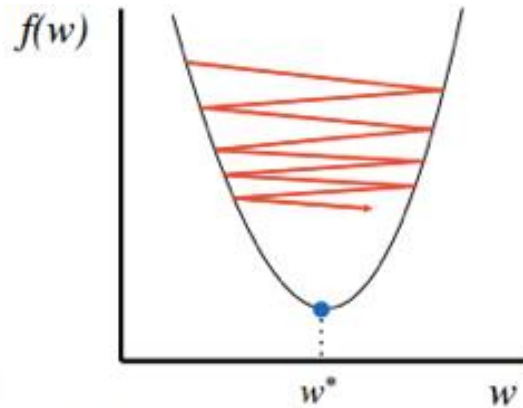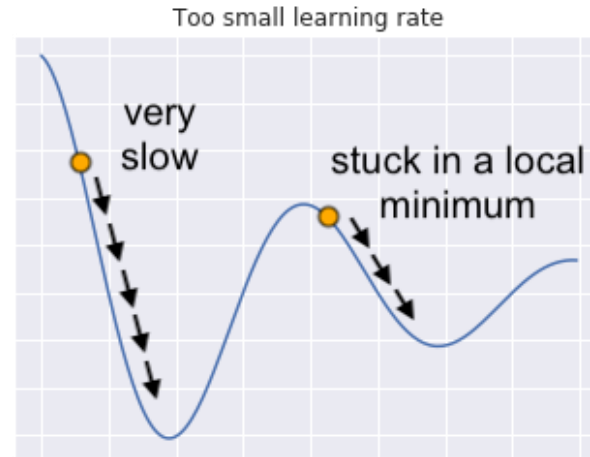**Graph Source**: https://cs231n.github.io/neural-networks-3/

# Learning Rate (cont'd)



$f(w)$

Too small: converge very slowly

Too big: overshoot and even diverge

Too small learning rate

very slow

stuck in a local minimum

Too large learning rate

divergence

**Graph Source**: https://towardsdatascience.com/the-learning-rate-finder-6618dfcb2025

**Graph Source**: https://srdas.github.io/DLBook/GradientDescentTechniques.html

One technique that can help the network out of local minima is the use of a **momentum** term.

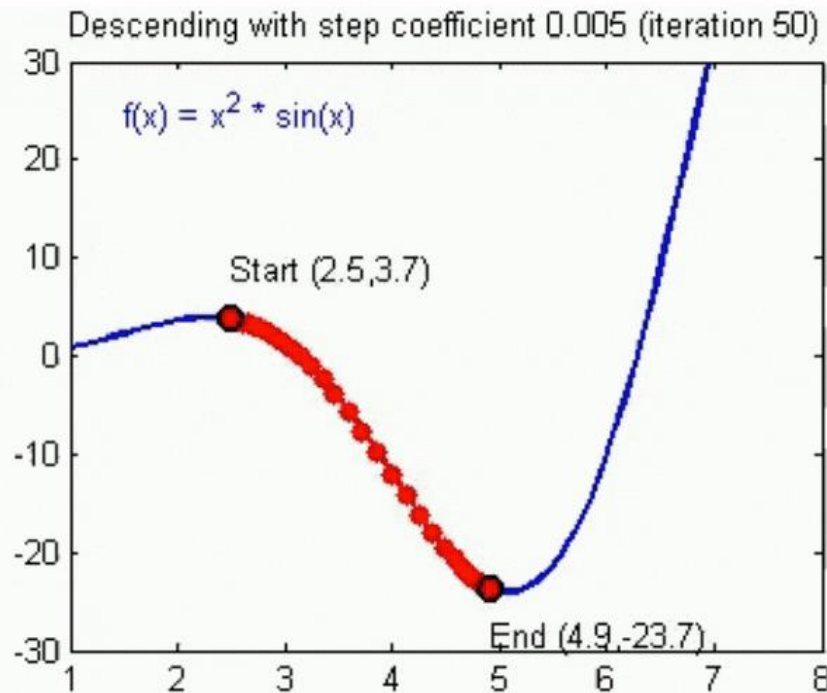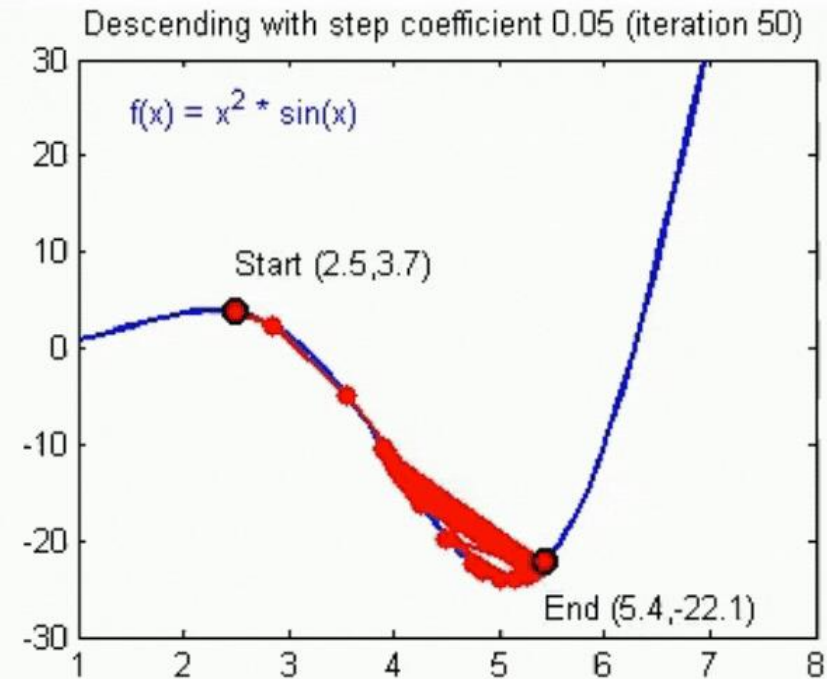$$\Delta w_{kj}(n) = \eta * \delta_k(n) * x_j(n) + \alpha \Delta w_{kj}(n-1)$$

Momentum factor

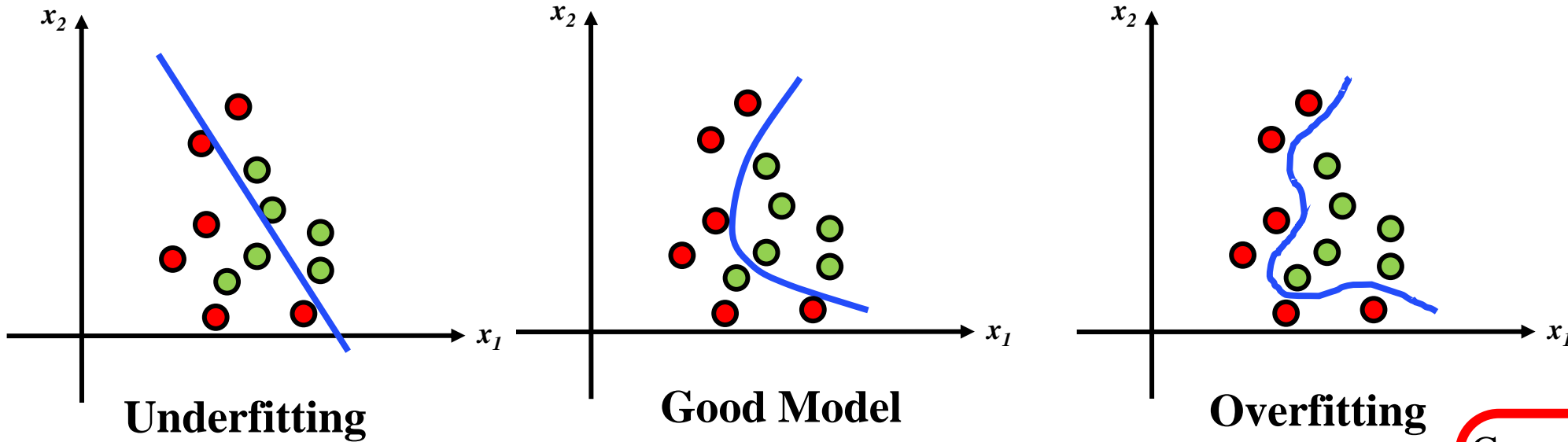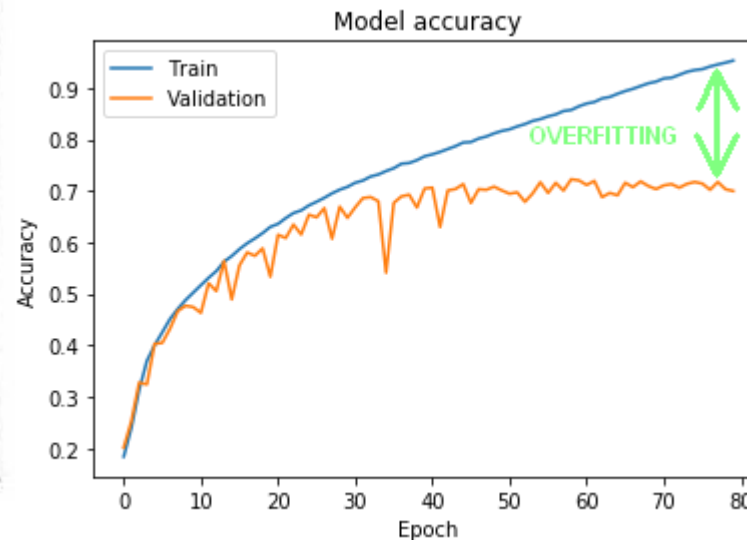Weight increment from previous iteration

46

**Convergence**

**Divergence**

# Overfitting



**Underfitting**

**Good Model**

**Overfitting**

**Can NOT be** **generalized**

Correctly classify test patterns it has never seen (learned) before when tested in real-world problem

**Solution**
➤ Early stopping
➤ Regularization (Dropout)

48

# Vanishing Gradient

❑ Deeper neural networks (i.e., with multiple hidden layers) are *difficult* to train (difficulty increases geometrically).

$$\delta_j = f'(\cdot) * \sum_{k=1}^{D} \delta_k * w_k \quad \text{[local gradient] x [upstream gradient]}$$

➢ The gradients get **smaller** and smaller when *backpropagating* the error.
➢ After few layers of propagation, the gradient disappears (***vanishes***)
➢ The parameters in the deep layer will be **almost static**

❑ Solution
➢ **Modify** the activation function
➢ **Use** batch normalization (sort of regularization)

# ANN Advantages and Disadvantages

❑ Advantages

➢ Very **simple** principles

➢ Highly parallel: *information processing is much more like the brain than a serial computer*

➢ Adapt to unknown situations, can model *complex* functions

➢ Ease of use, *learns by example*, and very little user domain-specific expertise needed.

❑ Disadvantages

➢ Very **complex** behaviors
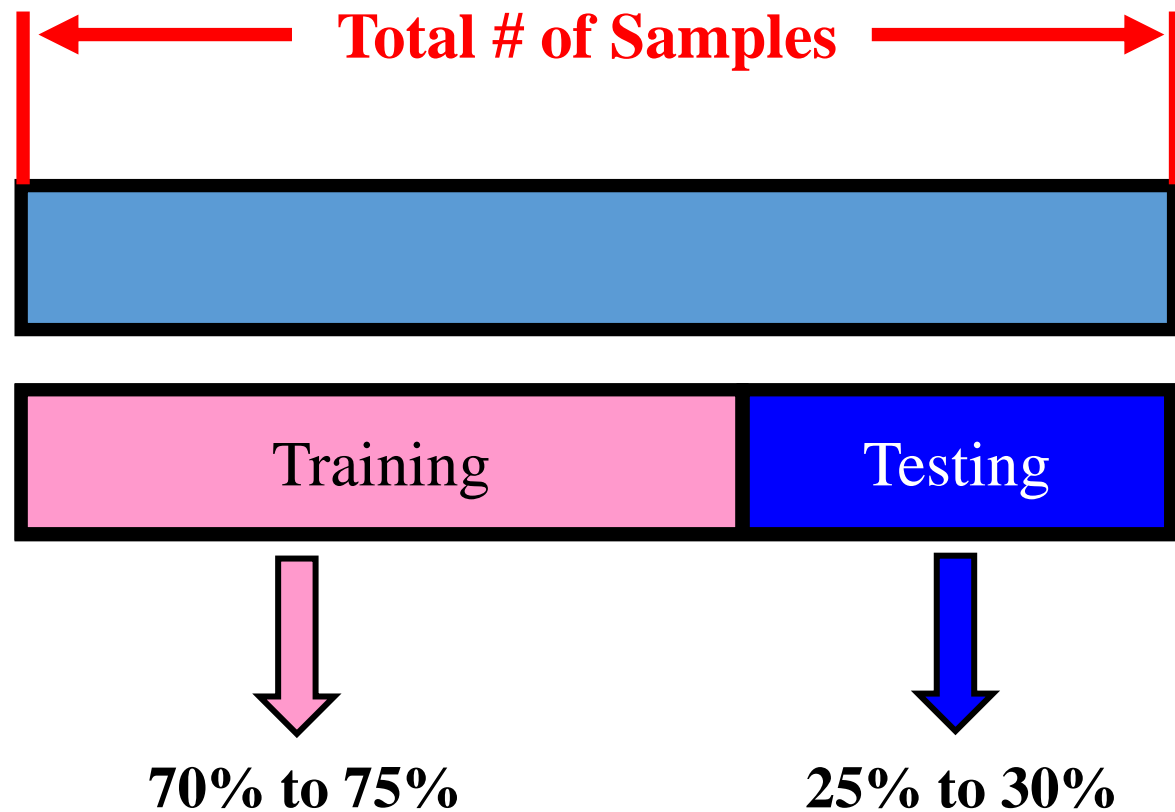
➢ Not exact.

➢ Needs training.

# ANN Terminology

❑ Neuron, unit (node)

❑ Weight and bias

❑ Transfer function (linear, sigmoid, ReLU, etc)

❑ Loss function (*mean squared error*, cross entropy, etc.)

❑ Learning rate, epoch, batch

❑ Backpropagation (***error*** propagation)

❑ Optimization (gradient descent (GD), stochastic GD, Adam,….etc.)

❑ Overfitting

❑ Dropout, Batch normalization
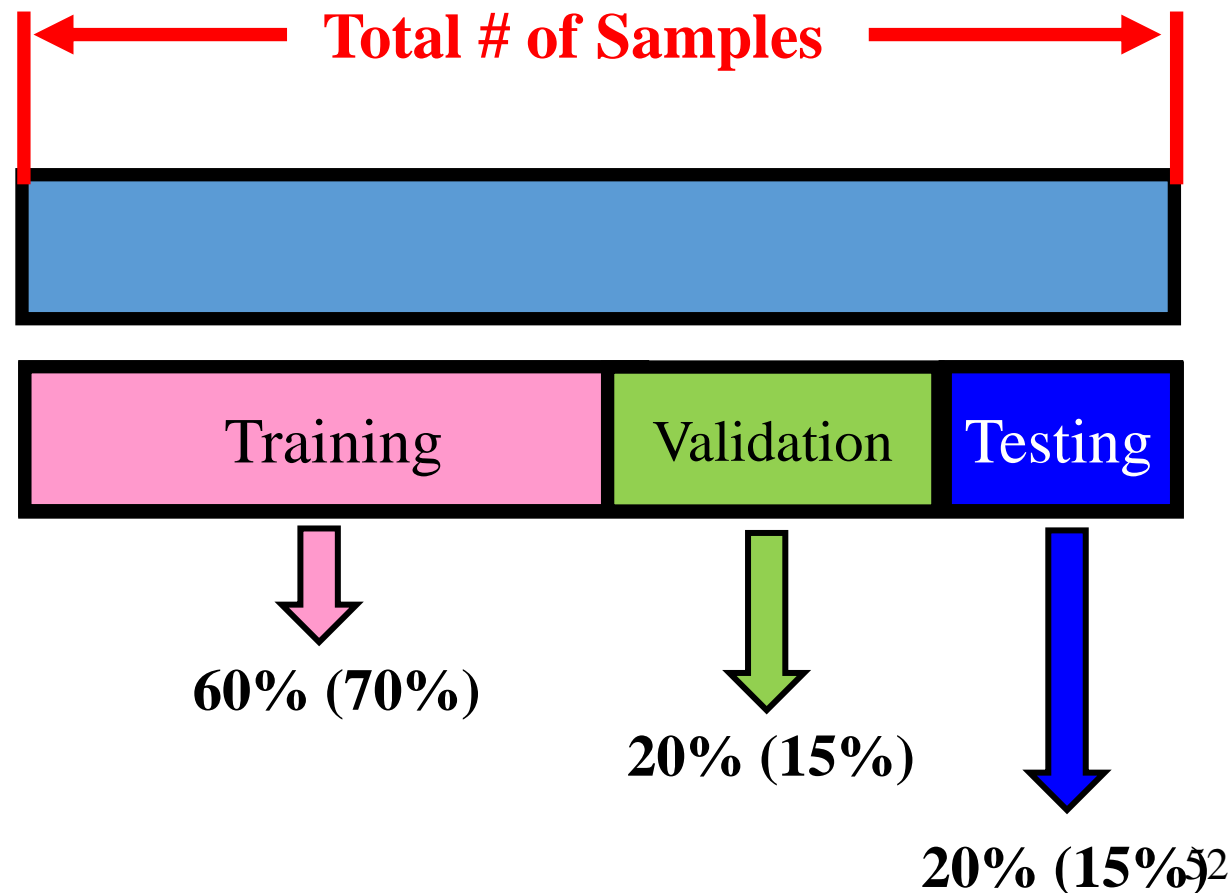
Each ANN aspect is considered a standalone research venue

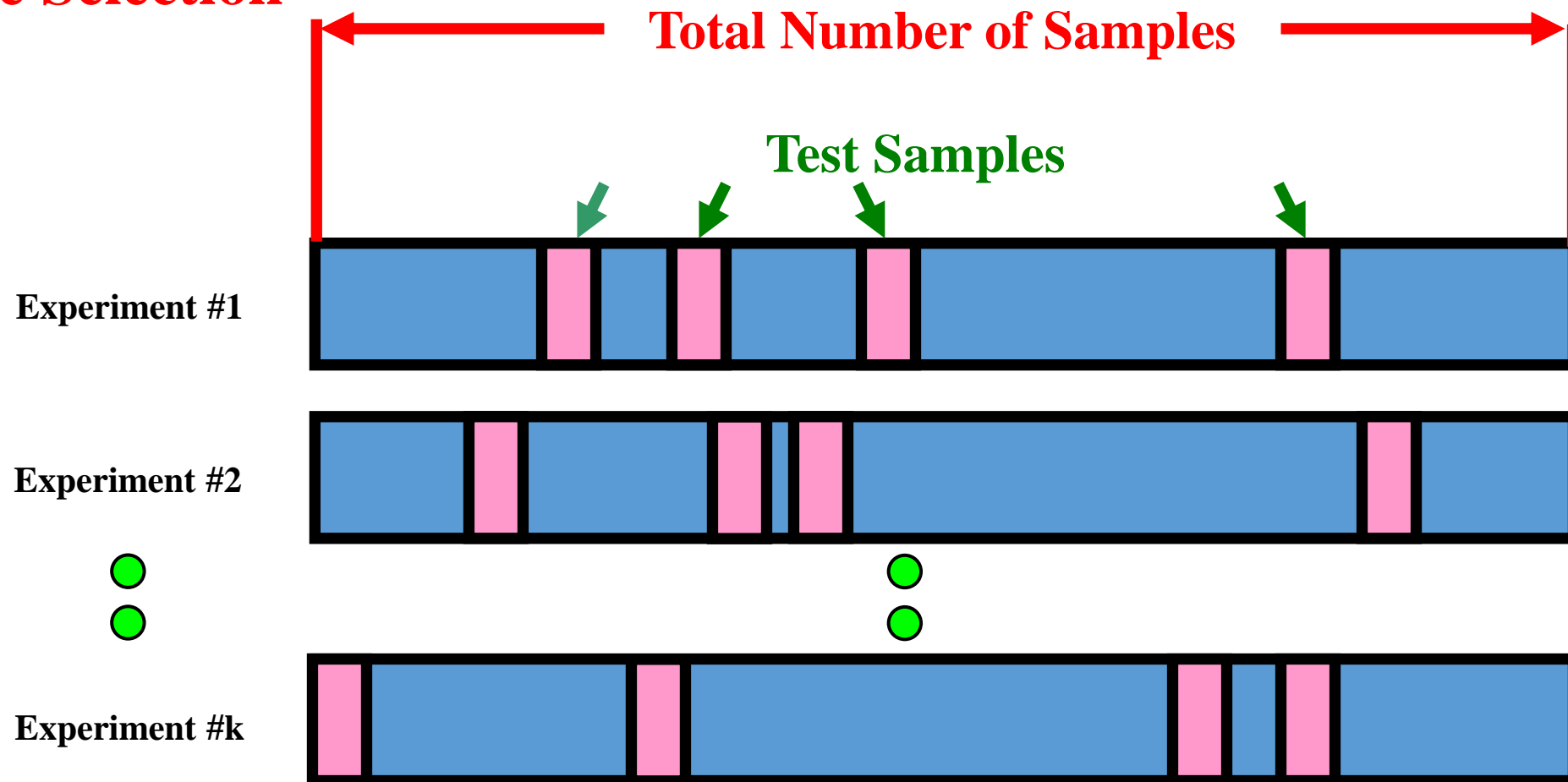# Validation Techniques

## Data Splitting

**Training/Testing**

**Training/Validation/Testing**

Total # of Samples

Total # of Samples

Training — Testing

Training — Validation — Testing

70% to 75%

25% to 30%

60% (70%)

20% (15%)

20% (15%)

# Validation Techniques

**Random Sample Selection**
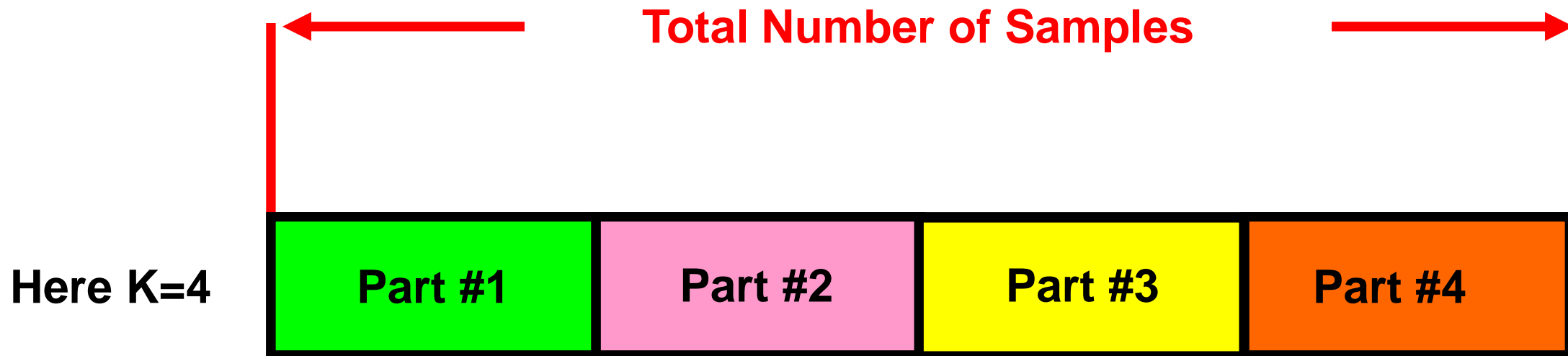


$$E = \frac{1}{k} \sum_{i=1}^{k} E_i$$

→ k is the number of experiment

→ $E_i$ is the average error for each experiment using only testing data
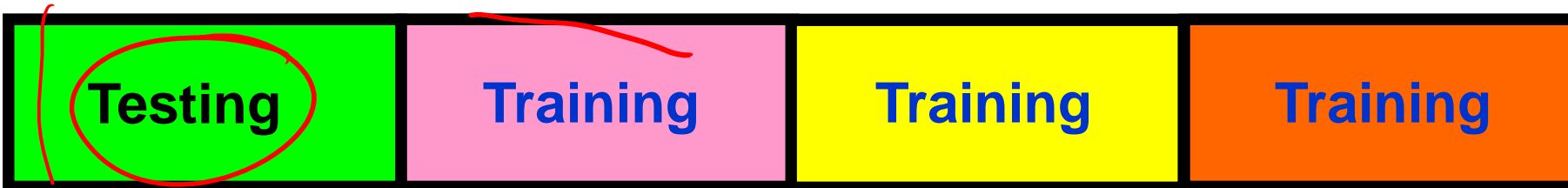
53

# Validation Techniques

## Cross Validation

Divide data into mutually exclusive and equal-sized subsets, folds, and this number is called K
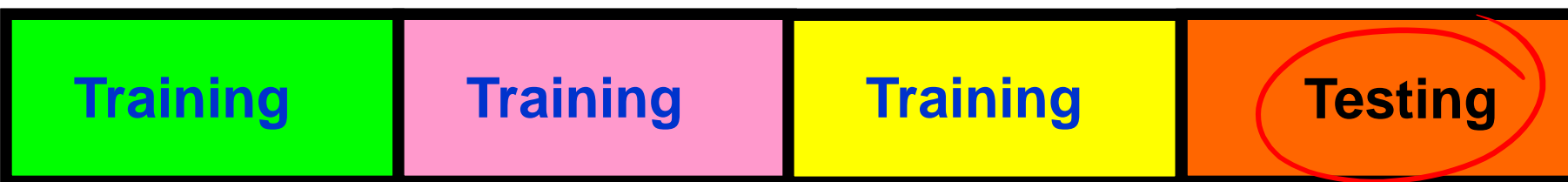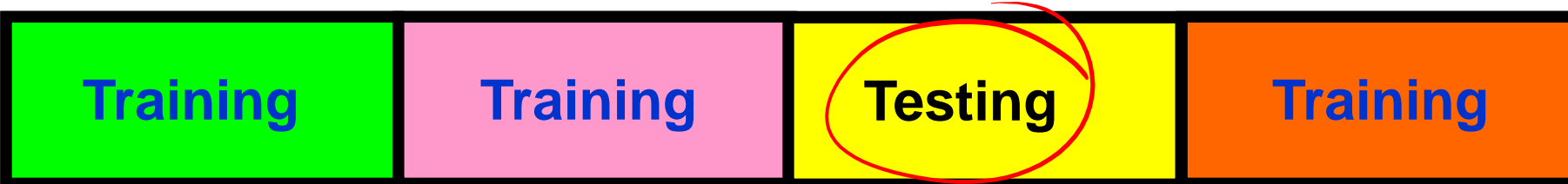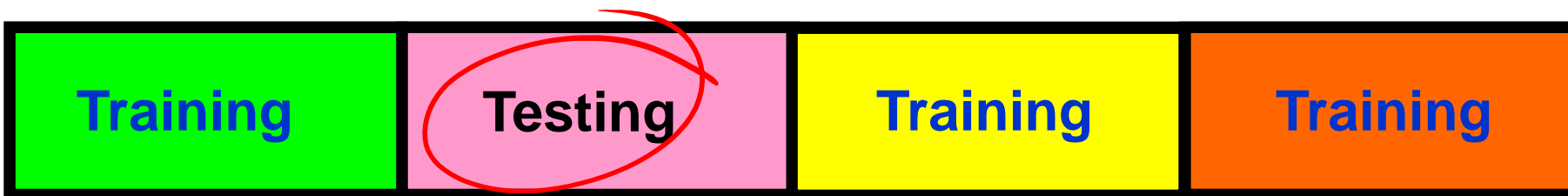
**Total Number of Samples**

**Here K=4**

| Part #1 | Part #2 | Part #3 | Part #4 |
|---------|---------|---------|---------|

$$E = \frac{1}{k} \sum_{i=1}^{k} E_i$$

→ k is the number of folds

→ $E_i$ is the average error for each fold

54

# Validation Techniques

75, 25

| Testing | Training | Training | Training |
|---------|----------|----------|----------|

| Training | Testing | Training | Training |
|----------|---------|----------|----------|

| Training | Training | Testing | Training |
|----------|----------|---------|----------|

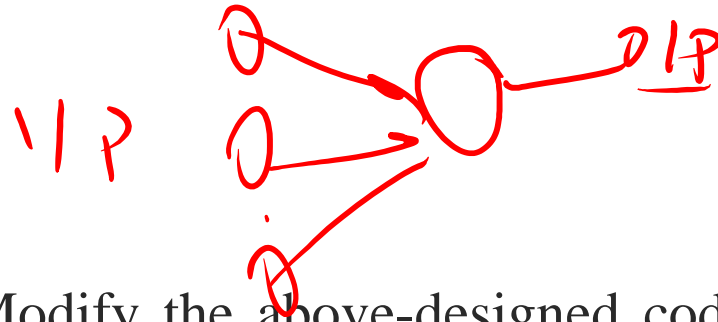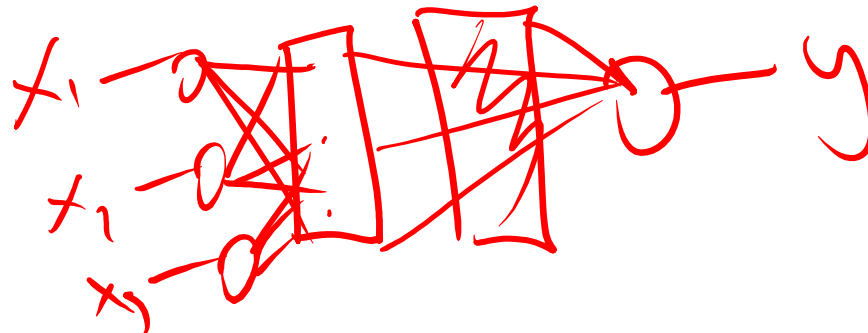| Training | Training | Training | Testing |
|----------|----------|----------|---------|

55

# Assignments

❑ **Assignment 1**: Design your own simple ANN, (one perceptron with one input layer and one output neuron). Use the data points listed in the adjacent Table as your training data. Assume the activation function is sigmoid and assume there is no bias for simplicity (b=0). Test your design using different iteration numbers.

| $X_1$ | $X_2$ | $X_2$ | d |
|-------|-------|-------|---|
| 0 | 0 | 1 | 0 |
| 0 | 1 | 1 | 1 |
| 1 | 0 | 1 | 1 |
| 1 | 1 | 1 | 0 |

❑ **Assignment 2**: Modify the above-designed code to implement a multi-layer perceptron, MLP (an ANN with one input layer, one hidden layer and one output layer) for the same data points above. Assume sigmoid activation function and there is no bias for simplicity (b=0). Test your approach using different iteration numbers and different number of nodes for the hidden layer (e.g., 4, 8, and 16).

❏ **Assignment 3**: Use the Keras library (*tensorflow.keras*) to build different ANNs using different numbers of hidden layers (**shallow**: 1 hidden, output layer, **deeper**: two hidden layers with 12 and 8 nodes respectively, and **more deep**: three hidden layers with 32, 16, 8 nodes respectively). Use the provided diabetic data sets (here) to train and test your design. Use the ReLU activation for the hidden layers and the sigmoid activation for the output neuron, loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'], epochs = 150.

❏ **Assignment 4**: Redo assignment #3 using 80% of the data for training and 20% of the data for testing. Also, plot the training accuracy and loss curves for your designed networks

# Thank You
# &
# Questions