# Artificial Intelligence

Instructor: Kietikul Jearanaitanakij
Department of Computer Engineering
King Mongkut's Institute of Technology Ladkrabang
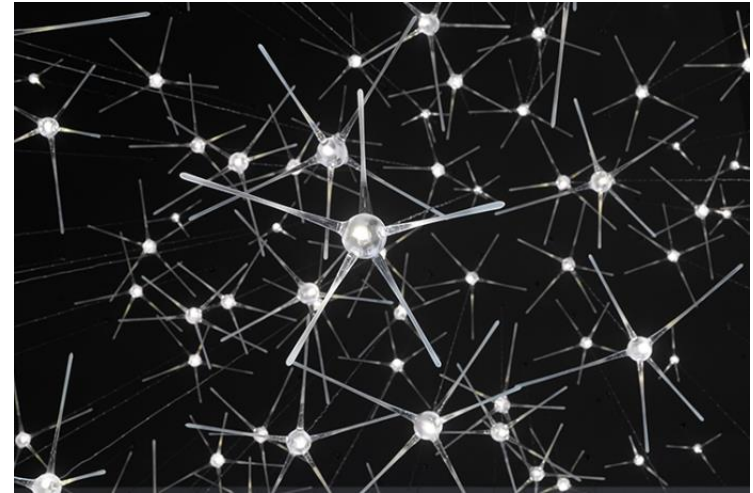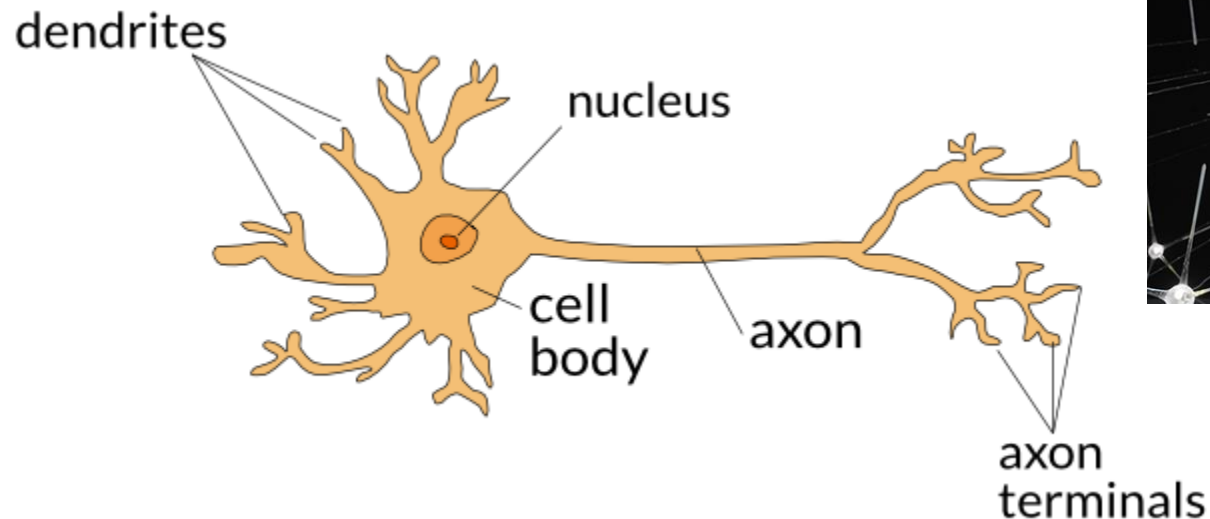
# L e c t u r e 9
## Perceptron And Artificial Neural Networks

- Linear Classification With Perceptron

- Perceptron learning algorithm

- Multilayer neural networks

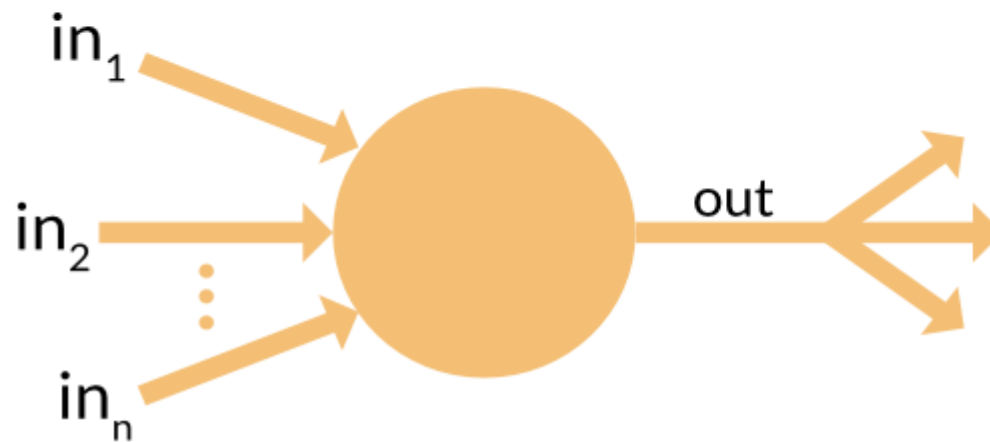- Backpropagation for multilayer NN

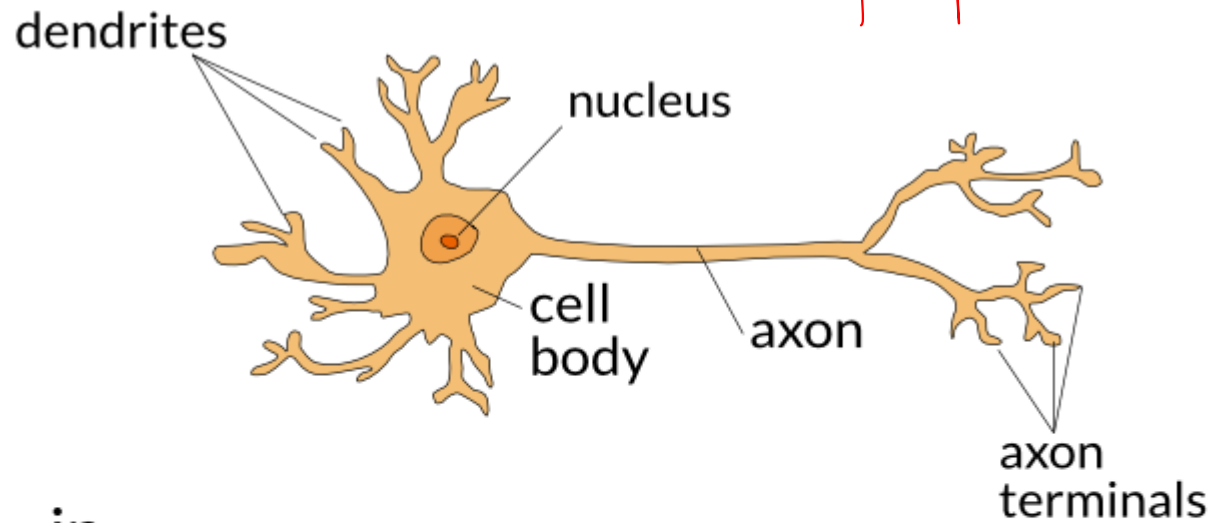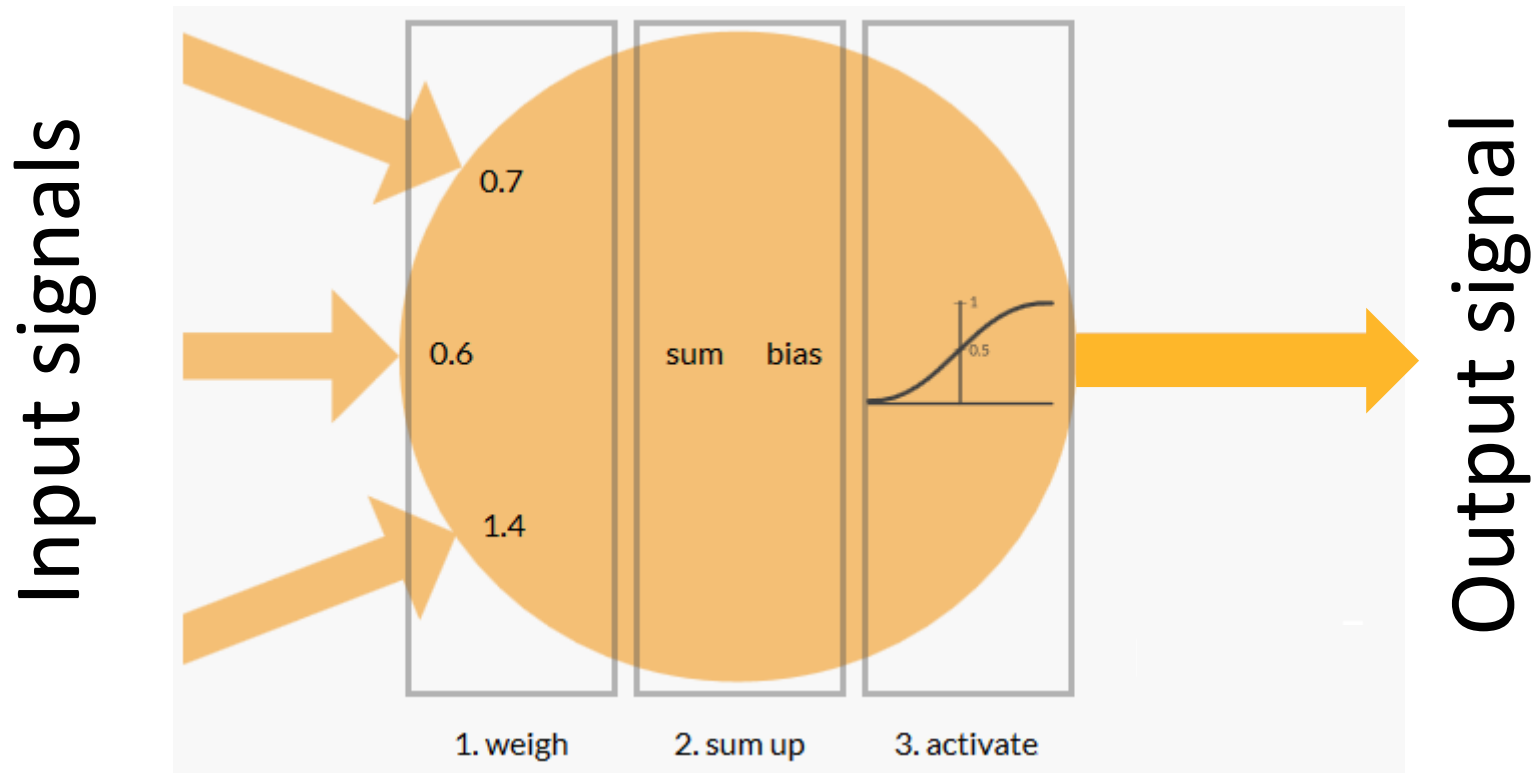# Linear Classification With Perceptron

**Biological neuron**



- Dendrites receive signals
- Axon sends signals out to other neurons

# Artificial neuron

perceptron ((i))



dendrites

nucleus

cell body

axon

axon terminals

$in_1$

$in_2$

$in_n$

out

# Inside artificial neuron



Input signals

Output signal

0.7

0.6

1.4

sum    bias

1. weigh    2. sum up    3. activate
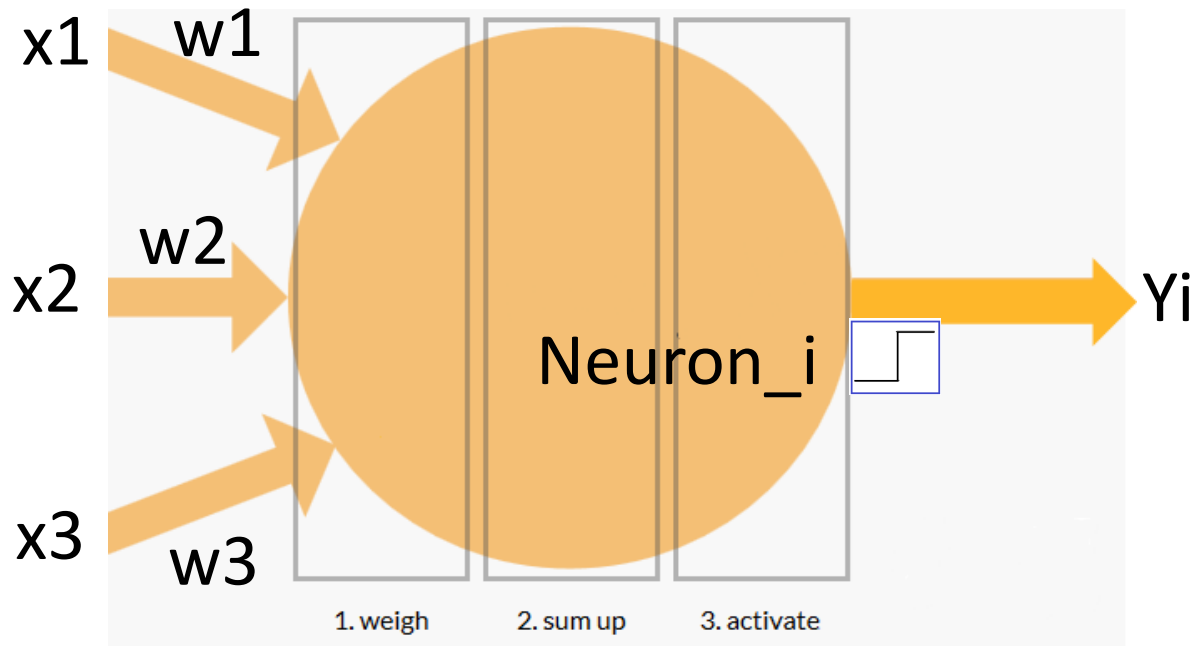
We will call this artificial neuron as a **perceptron**.

# Notations

**Perceptron** is the simplest form of a neural network. It consists of a single neuron with adjustable weights and a hard limit activation function.
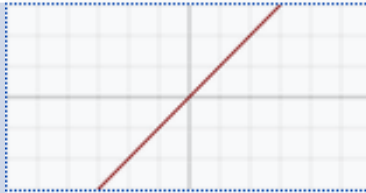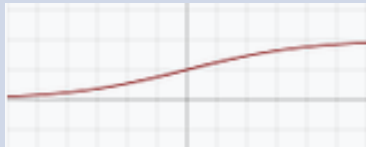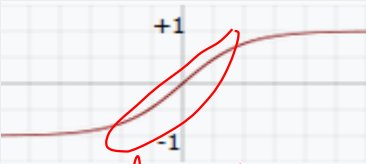


x1 w1

x2 w2

x3 w3

Neuron_i

Yi

1. weigh   2. sum up   3. activate

There are many kinds of activation function.

**Frank Rosenblatt**

invented perceptron in 1957

# Activation functions

| Name | Equation | Plot |
|------|----------|------|
| Identity (Linear) | $f(x) = x$  *useless* | |
| Binary step | $f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$ | *step* |
| Sigmoid (Logistic) | $f(x) = \dfrac{1}{1 + e^{-x}}$ | |
| TanH | $f(x) = \tanh(x) = \dfrac{2}{1 + e^{-2x}} - 1$ | +1 −1  *slope ไม่* |
| Softmax | $f_i(\vec{x}) = \dfrac{e^{x_i}}{\sum_{j=1}^{J} e^{x_j}}$ for $i = 1, \ldots, J$ | |

# Linear separability

Data in the n-dimensional space is **linearly separable** if two classes of data are divided by a hyperplane into two decision regions.

X2

Class +1

X1

Class 0

x1w1+x2w2-θ = 0

$ax + by + c$

X3

X1

X2

x1w1+x2w2+x3w3-θ = 0

In general, the hyperplane is defined by the linearly separable function.

$$\sum_{i=1}^{n} x_i w_i - \theta = 0$$

(The threshold θ can be used to shift the decision boundary.)

- The perceptron learns its classification task by making small adjustments in the weights to reduce the difference between the actual and desired (target) outputs of the perceptron.

Target value    Actual value

$$e(p) = Y_d(p) - Y(p)$$

Where   _error_   _desire_   _output ที่ออกมาจริง_

        p is the training pattern (example)

        $Y_d$(p) is the desired (target) output of pattern p

        Y(p) is the actual output

        e(p) is the difference (error) between $Y_d$(p) and Y(p)

_ค่าน้ำหนักใหม่ = ค่าน้ำหนักเก่า + ไปทางไหน x ระยะ-_
_( learning rate )_

- It uses e(p) to update weights of the next iteration,

$$w_i(p + 1) = w_i(p) + \alpha. x_i(p). e(p)$$

where α is the learning rate (0~1)

# Perceptron learning algorithm

**Step 1**: Initialization    เริ่มต้นให้ weight

– Set initial weights $w_1, w_2, \ldots, w_n$ and thresholds ($\theta$) to small random numbers in the range [-0.5,+0.5]

– p = 1   ; the first pattern    จำนวนแถบ ๆ

**Step 2**: Activation    ค่านอกเต้า

Step function:

$$Y(p) = step\left(\sum_{i=1}^{n} x_i(p).w_i(p) - \theta\right)$$

$$f(x) = \begin{cases} 0 & \text{for} \quad x < 0 \\ 1 & \text{for} \quad x \geq 0 \end{cases}$$

where n is the number of perceptron inputs.

**Step 3**: Weight training by using gradient descent

$$w_i(p + 1) = w_i(p) + \Delta w_i(p),$$

$$\Delta w_i(p) = \alpha.x_i(p).e(p) \quad \text{where } e(p) = Yd(p) - Y(p)$$
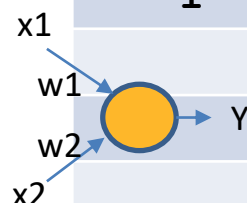
**Step 4**: Iteration  เริ่ม $p=2$

    Increase p by one, go back to step 2 until each pattern is trained.

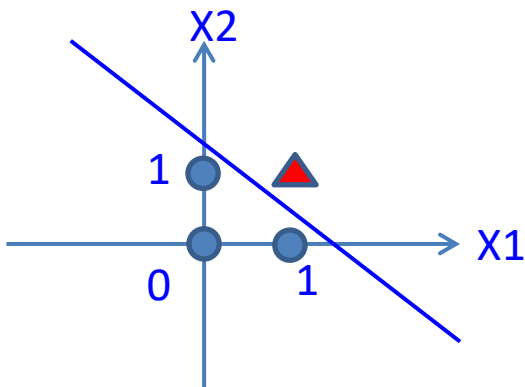**Step 5**: If the perceptron doesn't <u>converge</u>, p = 1 and repeat steps 2 – 4.

error

ลงแล้วนิ่ง → converge

epoch

e

# **Example**: Train a perceptron on AND

θ= 0.3, α = 0.1

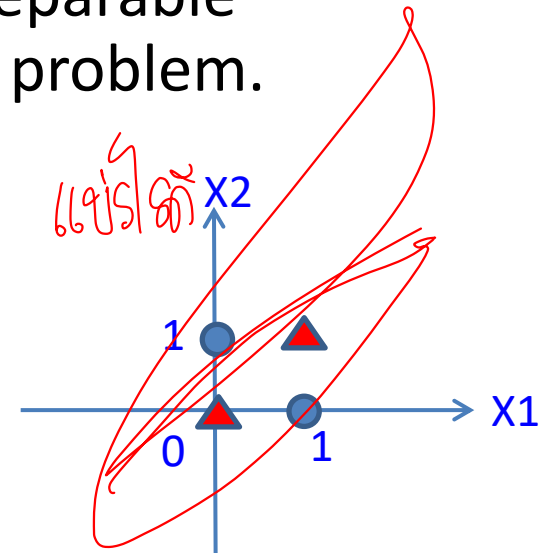| Epoch | Inputs | | Yd (desired) | Weights | | Y (actual) | Error | Weights | |
|---|---|---|---|---|---|---|---|---|---|
| | x1 | x2 | | w1 | w2 | | | w1 | w2 |
| 1 | 0 | 0 | 0 | 0.3 | −0.1 | 0 | 0 | 0.3 | −0.1 |
| | 0 | 1 | 0 | 0.3 | −0.1 | 0 | 0 | 0.3 | −0.1 |
| | 1 | 0 | 0 | 0.3 | −0.1 | 1 | −1 | 0.2 | −0.1 |
| | 1 | 1 | 1 | 0.2 | −0.1 | 0 | 1 | 0.3 | 0.0 |
| 2 | 0 | 0 | 0 | 0.3 | 0.0 | 0 | 0 | 0.3 | 0.0 |
| | 0 | 1 | 0 | 0.3 | 0.0 | 0 | 0 | 0.3 | 0.0 |
| | 1 | 0 | 0 | 0.3 | 0.0 | 1 | −1 | 0.2 | 0.0 |
| | 1 | 1 | 1 | 0.2 | 0.0 | 1 | 0 | 0.2 | 0.0 |
| . | . | | . | . | | . | . | . | |
| . | . | | . | . | | . | . | . | |
| . | . | | . | . | | . | . | . | |
| 5 | 0 | 0 | 0 | 0.1 | 0.1 | 0 | 0 | 0.1 | 0.1 |
| | 0 | 1 | 0 | 0.1 | 0.1 | 0 | 0 | 0.1 | 0.1 |
| | 1 | 0 | 0 | 0.1 | 0.1 | 0 | 0 | 0.1 | 0.1 |
| | 1 | 1 | 1 | 0.1 | 0.1 | 1 | 0 | 0.1 | 0.1 |

x1
w1
Y
w2
x2

# Problem of perceptron

- Perceptron can learn only simple linear separable problems, e.g., AND, OR. It failed on XOR problem.
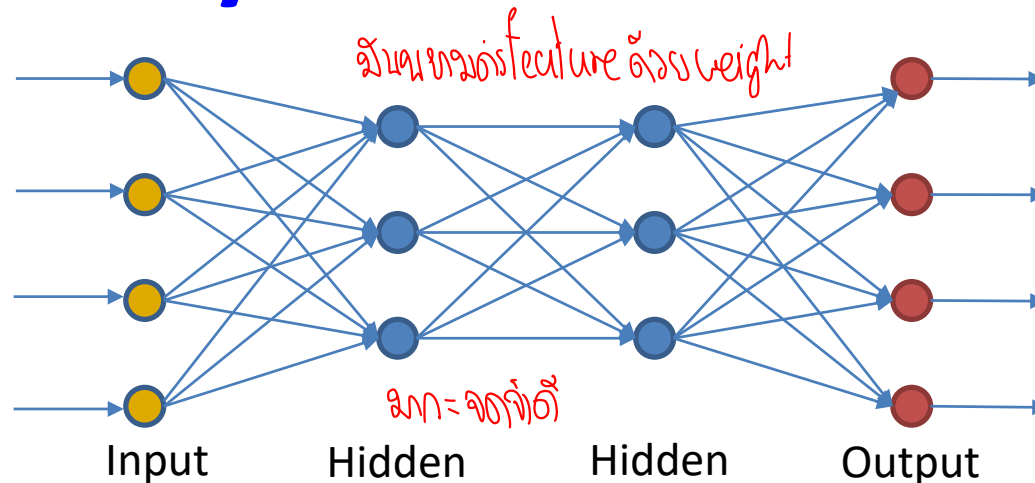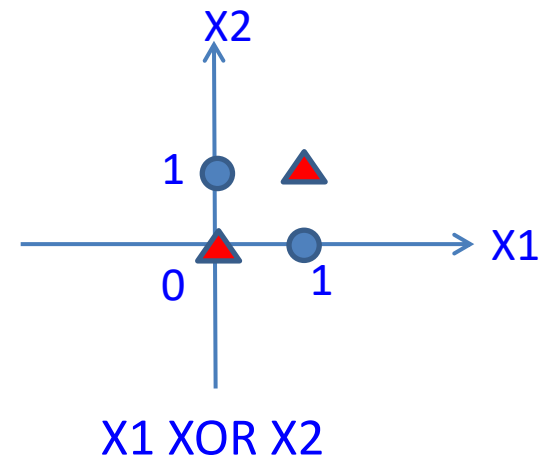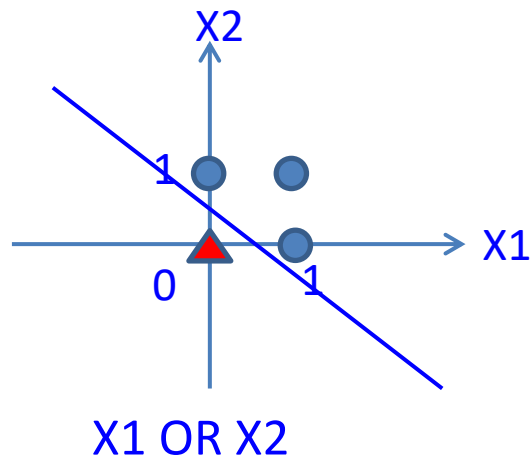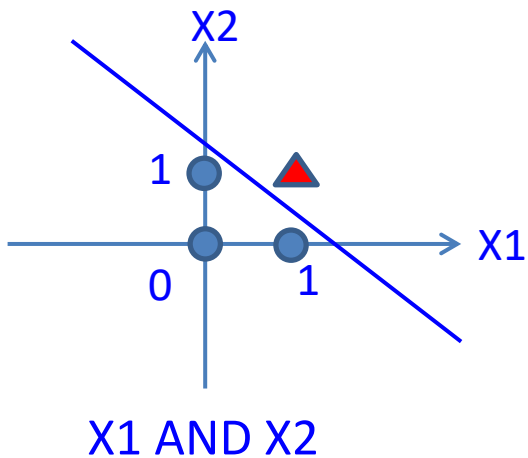


X1 AND X2

X1 OR X2

X1 XOR X2

- A single perceptron can classify only linear separable problems, regardless of whether we use a hard-limit or soft-limit activation functions.

- Moreover, increasing the number of perceptrons in the same layer doesn't help.

13

# Multilayer neural networks



- **Input layer** accepts input signals from the outside world and redistributes these signals to all neurons in the hidden layer.

- **Hidden layers** detect the features of the input patterns by adjusting weights of the neurons.

- **Output layer** accepts signals from the hidden layer and establishes the output pattern of the network.

- Multilayer neural networks can solve the **non-linearly separable problem**.

14

# Nonlinearly separable problem

X2

1 ●  ▲

0 ●  ● 1  → X1

X1 AND X2

X2

1 ●  ●

0 ▲  ● 1  → X1

X1 OR X2

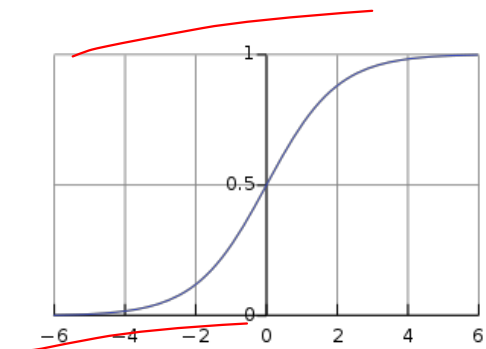X2

1 ●  ▲

0 ▲  ● 1  → X1

X1 XOR X2

Linearly separable

Non-linearly separable

- With one hidden layer, we can represent the continuous and simple discontinuous functions.
- With two hidden layers, even discontinuous function can be represented. จับแปะ

  ปกติไม่มีใครรอดเกิน 3 ถ้ายังไม่จบ
- Most practical application use three-layer neural network (1-1-1: input-hidden-output) to learn patterns. เขาจะก่อนๆ เป็น hidden เยอะๆ
- The most popular learning algorithm is "backpropagation" (Bryson & Ho, 1969).
- Each neuron determines its output Y as the following:

การถ่วงน้ำหนัก 1 neuron

Threshold (Interception)

$$X = \sum_{i=1}^{n} x_i w_i - \theta,$$

$$y = \frac{1}{1+e^{-x}} \quad ; 0 < y < 1 \quad \text{(Sigmoid)}$$



16

# Notations

$$X = \sum_{i=1}^{n} x_i w_i - \theta,$$

# **Backpropagation for multilayer NN**

- To propagate error signals, we start at the output layer and work backward to the hidden layer.

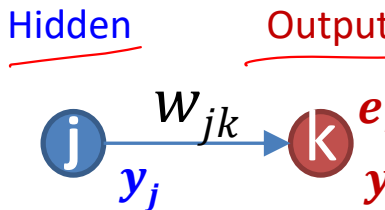- The error signal at the output of neuron k at iteration p is defined by:

$$e_k(p) = yd_k(p) - y_k(p)$$

    Where $yd_k(p)$ is the desired output of neuron k at iteration p.

- To update weights, the weight correction ($\Delta w$) is adjusted to the previous weights.
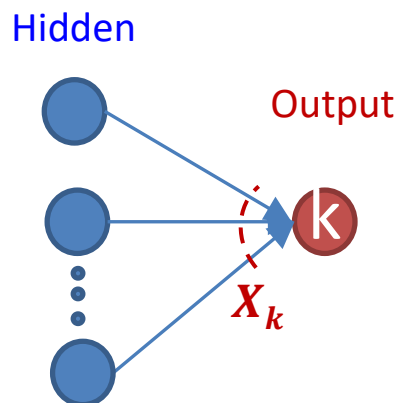
$$w_{jk}(p + 1) = w_{jk}(p) +\ \Delta w_{jk}(p),$$

**Hidden**   **Output**

j $\xrightarrow{w_{jk}}$ k   $e_k = yd_k - y_k$

$y_j$   $y_k$

$$\Delta w_{jk}(p) =\ \alpha . y_j(p). \delta_k(p),$$

*input*

*ปรับ weight base*

$$\delta_k(p) = \frac{\partial y_k(p)}{\partial X_k(p)}. e_k(p),$$

Error gradient at neuron k

*sigmoid × error*

*gradient ความชัน → dlt*

*ชนใกล้ error*

*เอียง*

For a sigmoid activation function ($y_k =\ \frac{1}{1+e^{-X_k}}$),

*ทดแทน (sigmoid)*

$$\frac{\partial y_k(p)}{\partial X_k(p)} = \frac{\partial\left[\frac{1}{1+e^{-X_k(p)}}\right]}{\partial X_k(p)} = \frac{e^{-X_k(p)}}{\left(1+e^{-X_k(p)}\right)^2}$$

$$\delta_k(p) = y_k(p). \left(1 - y_k(p)\right). e_k(p)$$

*sigmoid (1 − sigmoid) · error*

**Hidden**

**Output**

k

$X_k$

- To update weights in the hidden layer, we do similar way:

$$w_{ij}(p+1) = w_{ij}(p) + \Delta w_{ij}(p),$$

$$\Delta w_{ij}(p) = \alpha . x_i(p) . \delta_j(p),$$

$$\delta_j(p) = y_j(p) . \left(1 - y_j(p)\right) . \sum_{q=1}^{l} [\delta_{kq}(p) . w_{jkq}(p)]$$

$$\frac{\partial y_j(p)}{\partial X_j(p)}$$

Where $l$ is the number of neurons in the output layer.

$$y_j(p) = \frac{1}{1+e^{-X_j(p)}},$$

$$X_j(p) = \sum_{i=1}^{n} [x_i(p) . w_{ij}(p)] - \theta_j$$

# Complete training algorithm

**Step1:** Weights Initialization
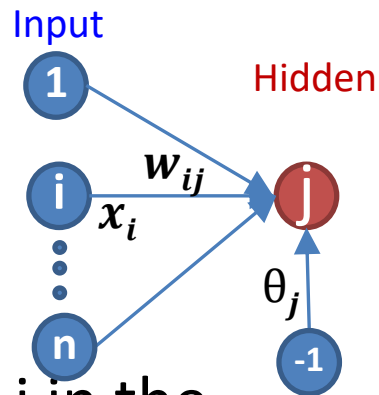
 - Set all weights and thresholds (θ) to random number uniformly distributed inside a small range, e.g., (-0.5, +0.5).

**Step2:** Activation (Feed forward computation)

 - Calculate the actual outputs of neurons in the hidden layer.

$$y_j(p) = sigmoid \left[ \sum_{i=1}^{n} (x_{i(p)} . w_{ij(p)}) - \theta_j \right]$$

 Where n is the number of inputs of neuron j in the hidden layer.

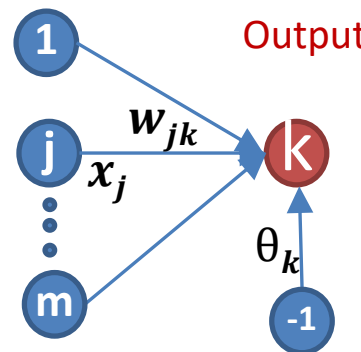- Calculate the actual outputs of neurons in the output layer.

$$y_k(p) = sigmoid\left[\sum_{j=1}^{m}\left(x_{j(p)} \cdot w_{jk(p)}\right) - \theta_k\right]$$

Where m is the number of inputs of neuron k in the output layer.

# **Step3:** Weight training (Backpropagation)

- Calculate the error gradient of neurons in the hidden-output layer.

$$e_k(p) = yd_k(p) - y_k(p), \text{ ค่า error}$$

$$\delta_k(p) = y_k(p).\big(1 - y_k(p)\big).e_k(p), \text{ ค่าความชัน}$$

$$\Delta w_{jk}(p) = \alpha.y_j(p).\delta_k(p), \text{ ตัวที่จะขยับ}$$

$$w_{jk}(p+1) = w_{jk}(p) + \Delta w_{jk}(p) \quad \text{เริ่มขยับ} \leftarrow \text{อัปเดตเวปา-ดิม}$$

$$\text{backprop ของ state ย้อน}$$

**Hidden**   **Output**

j  $y_j$  $w_{jk}$  →  k   $e_k = yd_k - y_k$

$y_k$

- Calculate the error gradient of neurons in the  input-hidden layer.

$$\delta_j(p) = y_j(p).\big(1 - y_j(p)\big).\sum_{k=1}^{l} \delta_k(p).w_{jk}(p), \text{ การขยับรวม}$$

$$\Delta w_{ij}(p) = \alpha.x_i(p).\delta_j(p), \text{ ขยับ}$$

$$w_{ij}(p+1) = w_{ij}(p) + \Delta w_{ij}(p)$$

เริ่มขยาย

**Input**   **Hidden**   **Output**

$x_i$  $w_{ij}$  →  j  $\delta_j$  $y_j$

$w_{jk1}$ → k1  $\delta_{k1}$

$w_{jk2}$ → k2  $\delta_{k2}$

$w_{jkl}$ → kl  $\delta_{kl}$

จะออกรับก่อนเปิด
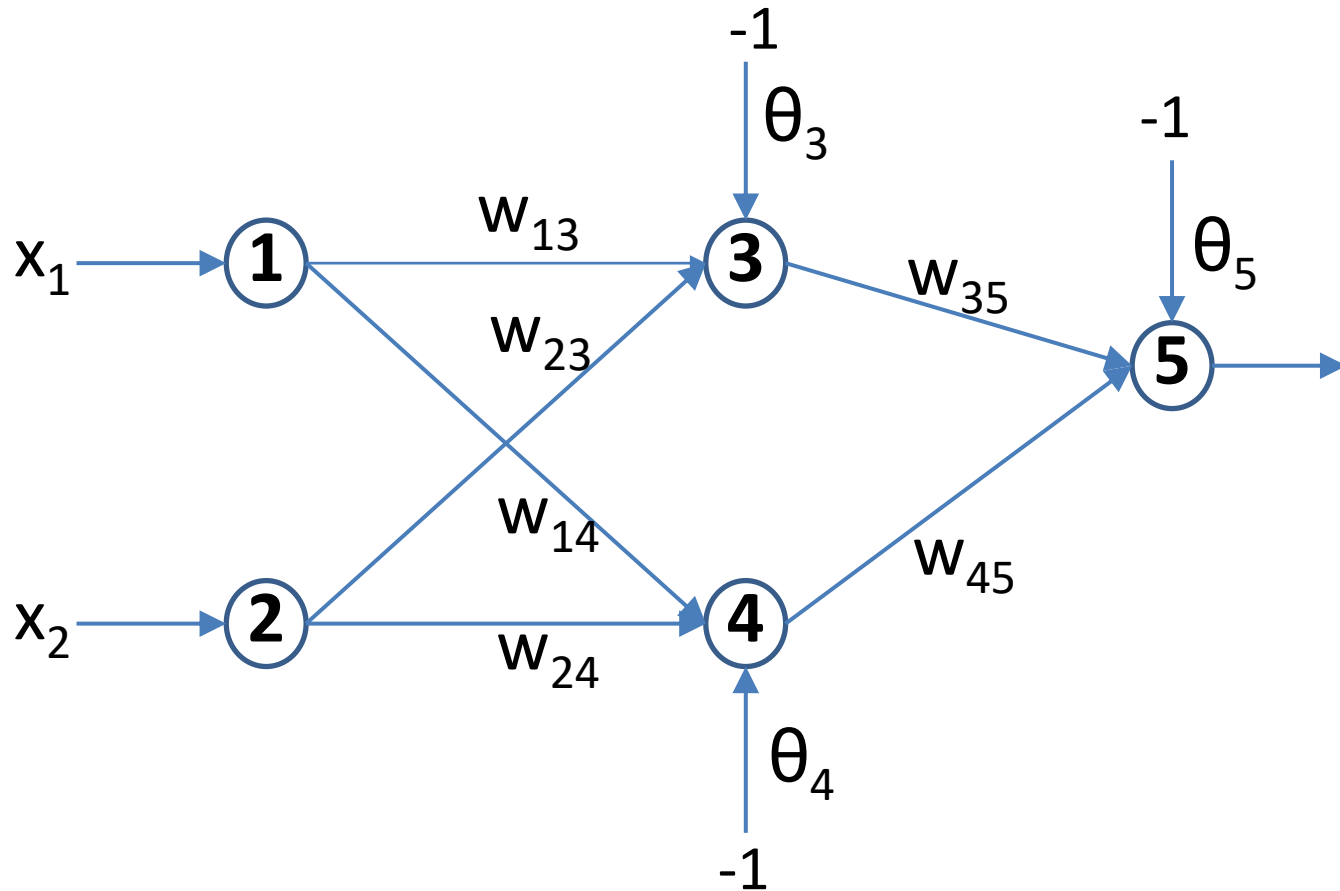
# Step4: Iterations

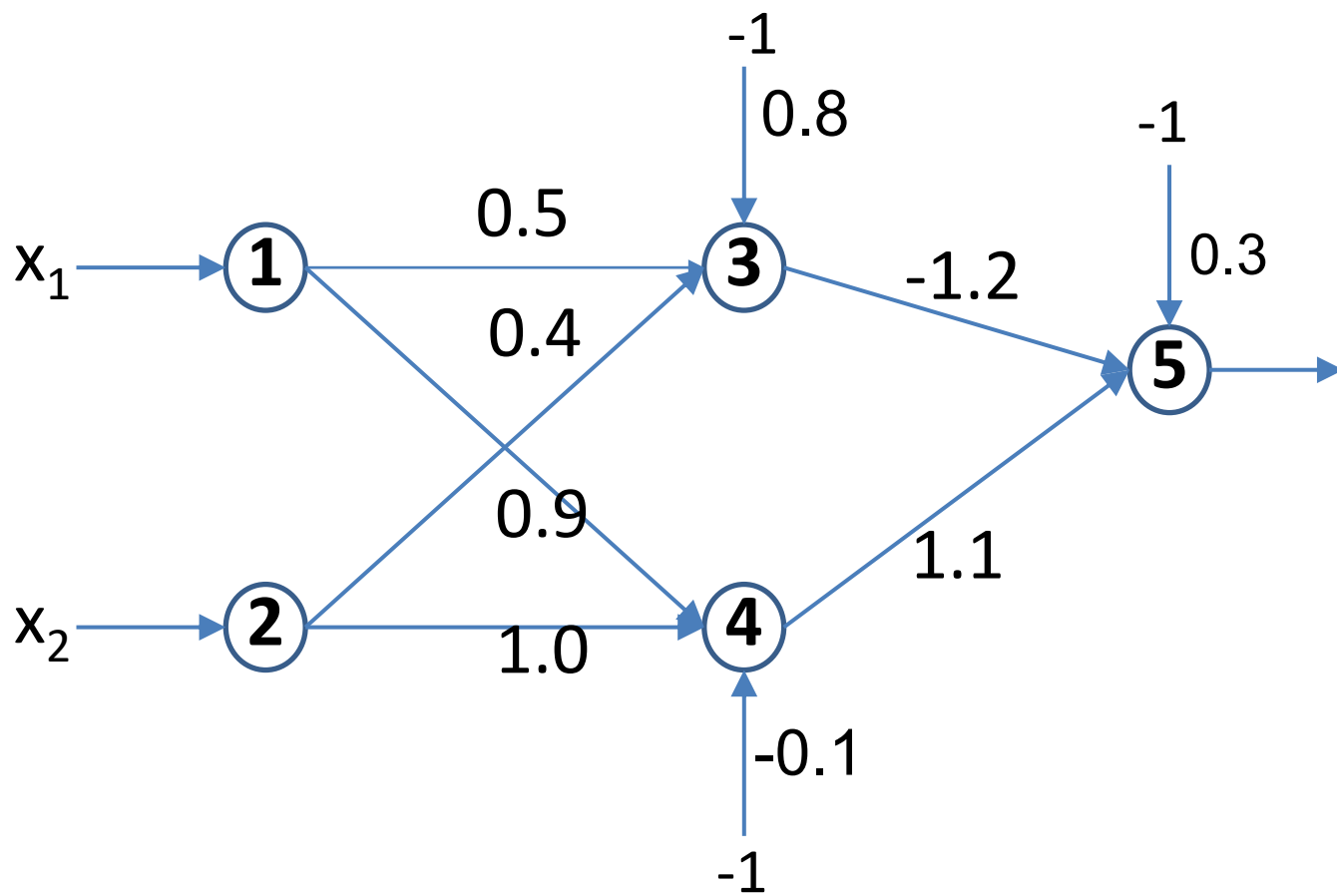    - Take the next training pattern, p+1, and go back to step 2. Then repeat the process until the selected error criterion is satisfy.

    - A simple stop criterion is when the sum-squared error (SSE) less than a certain number, e.g., 0.1.

$$SSE = \sum_{p=1}^{\#patterns} \sum_{k=1}^{\#outputs} [yd_k(p) - y_k(p)]^2$$
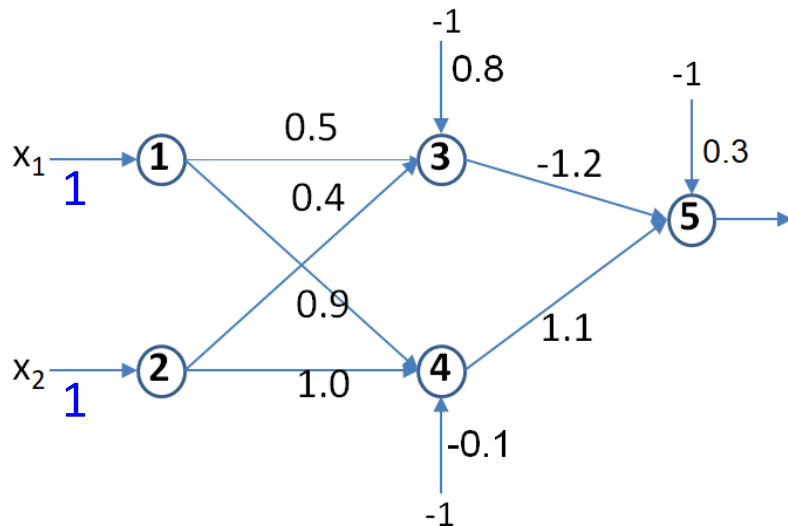
# Example: 3-layer NN for XOR problem



Recall that a single-layer perceptron cannot solve XOR problem.

All weights and thresholds are randomly initialized as follows:

| Parameter | Initial Value |
| --- | --- |
| $w_{13}$ | 0.5 |
| $w_{14}$ | 0.9 |
| $w_{23}$ | 0.4 |
| $w_{24}$ | 1.0 |
| $w_{35}$ | -1.2 |
| $w_{45}$ | 1.1 |
| $\theta_3$ | 0.8 |
| $\theta_4$ | -0.1 |
| $\theta_5$ | 0.3 |

- Consider the training pattern, <1, 1, 0>



$x_1$ 1

$x_2$ 1

$$y_3 = \frac{1}{[1 + e^{-(1*0.5+1*0.4-1*0.8)}]} = 0.5250$$

ตรวจดูไป

$$y_4 = \frac{1}{[1 + e^{-(1*0.9+1*1.0+1*0.1)}]} = 0.8808$$

$$y_5 = \frac{1}{[1 + e^{-(-0.525*1.2+0.8808*1.1-1*0.3)}]} = 0.5097$$

ค่าที่ต้องการ      ค่าที่ได้จริง

$$\text{e} = yd_5 - y_5 = 0 - 0.5097 = -0.5097$$

$$\delta_5 = y_5 \cdot (1 - y_5) \cdot e$$
$$= 0.5097 * (1 - 0.5097) * (-0.5097)$$
$$= -0.1274$$

28

- Assume that the learning rate, $\alpha$, is equal to 0.1.

$$\Delta w_{35} = \alpha . y_3 . \delta_5 = 0.1 * 0.5250 * (-0.1274) = -0.0067$$

$$\Delta w_{45} = \alpha . y_4 . \delta_5 = 0.1 * 0.8808 * (-0.1274) = -0.0112$$

$$\Delta \theta_5 = \alpha . (-1) . \delta_5 = 0.1 * (-1) * (-0.1274)$$

$$= 0.0127$$



- Next, we calculate the error gradient for neuron 3 and 4.

$$\delta_3 = y_3 . (1 - y_3) . \delta_5 . w_{35}$$
$$= 0.525 * (1 - 0.525) * (-0.1274) * (-1.2) = 0.0381$$

$$\delta_4 = y_4 . (1 - y_4) . \delta_5 . w_{45}$$
$$= 0.8808 * (1 - 0.8808) * (-0.1274) * 1.1 = -0.0147$$

- Calculating Δ of weights and thresholds



$$\Delta w_{13} = \alpha . x_1 . \delta_3 = 0.1 * 1 * 0.0381 = 0.0038$$

$$\Delta w_{23} = \alpha . x_2 . \delta_3 = 0.1 * 1 * 0.0381 = 0.0038$$

$$\Delta \theta_3 = \alpha . (-1) . \delta_3 = 0.1 * (-1) * 0.0381 = -0.0038$$

$$\Delta w_{14} = \alpha . x_1 . \delta_4 = 0.1 * 1 * (-0.0147) = -0.0015$$

$$\Delta w_{24} = \alpha . x_2 . \delta_4 = 0.1 * 1 * (-0.0147) = -0.0015$$

$$\Delta \theta_4 = \alpha . (-1) . \delta_4 = 0.1 * (-1) * (-0.0147) = 0.0015$$

- Lastly, we update all weights and thresholds in the network.

$$w_{13} = w_{13} + \Delta w_{13} = 0.5 + 0.0038 = 0.5038$$

$$w_{14} = w_{14} + \Delta w_{14} = 0.9 - 0.0015 = 0.8985$$

$$w_{23} = w_{23} + \Delta w_{23} = 0.4 + 0.0038 = 0.4038$$

$$w_{24} = w_{24} + \Delta w_{24} = 1.0 - 0.0015 = 0.9985$$

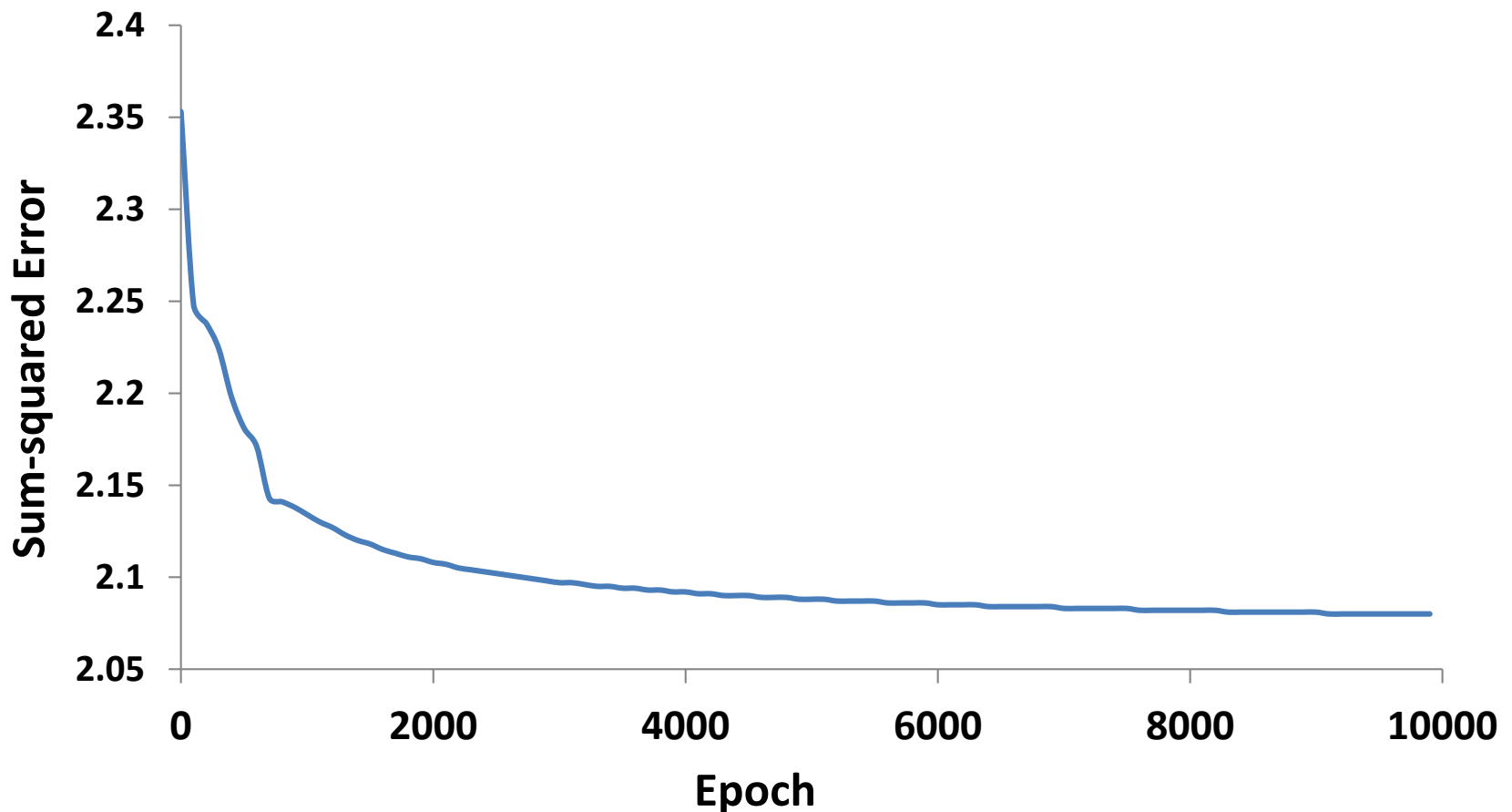$$w_{35} = w_{35} + \Delta w_{35} = -1.2 - 0.0067 = -1.2067$$

$$w_{45} = w_{45} + \Delta w_{45} = 1.1 - 0.0112 = 1.0888$$

$$\theta_3 = \theta_3 + \Delta\theta_3 = 0.8 - 0.0038 = 0.7962$$

$$\theta_4 = \theta_4 + \Delta\theta_4 = -0.1 + 0.0015 = -0.0985$$

$$\theta_5 = \theta_5 + \Delta\theta_5 = 0.3 + 0.0127 = 0.3127$$

- Repeat the same computation for all training patterns (1 epoch)
- Repeat the process for another epoch until the sum of squared error (SSE) is less than a certain number, e.g. 0.001.

- Sum of squared errors (SSE) of the final network

| Input | | Desired Output | Actual Output | Error | SSE |
|---|---|---|---|---|---|
| $x_1$ | $x_2$ | $y_d$ | $y_5$ | $e$ | |
| 1 | 1 | 0 | 0.0155 | -0.0155 | 0.0010 |
| 0 | 1 | 1 | 0.9849 | 0.0151 | |
| 1 | 0 | 1 | 0.9849 | 0.0151 | |
| 0 | 0 | 0 | 0.0175 | -0.0175 | |