

# Neural Network

## MATLAB進階程式語言與實作

盧家鋒 Chia-Feng Lu, Ph.D.  
Department of Biomedical Imaging and  
Radiological Sciences, NYCU  
[alvin4016@nycu.edu.tw](mailto:alvin4016@nycu.edu.tw)

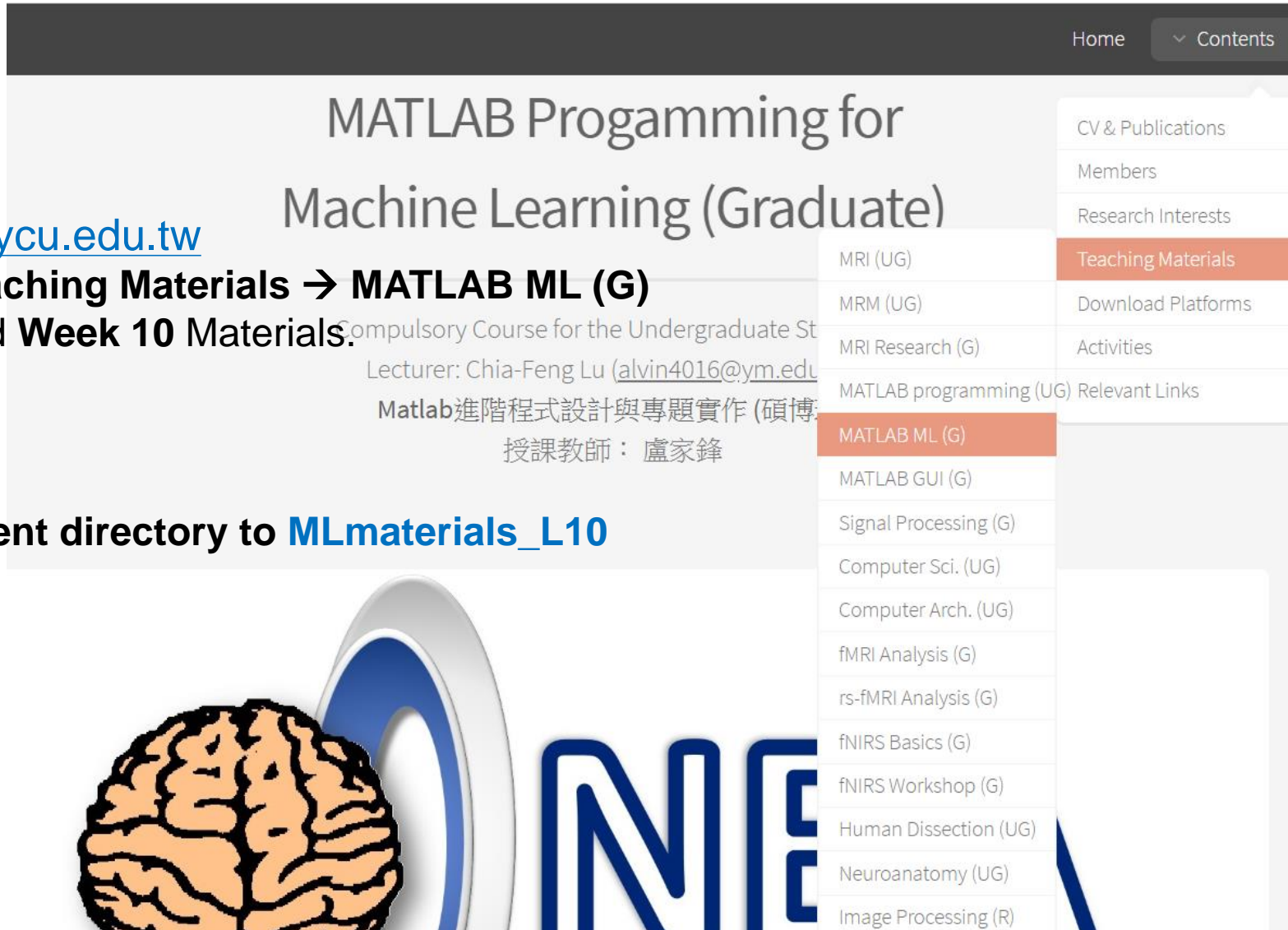
# Teaching Materials

<http://cflu.lab.nycu.edu.tw>

**Contents → Teaching Materials → MATLAB ML (G)**

Please download **Week 10** Materials.

Please set current directory to **MLmaterials\_L10**



Home Contents

## MATLAB Programming for Machine Learning (Graduate)

Compulsory Course for the Undergraduate Students

Lecturer: Chia-Feng Lu ([alvin4016@ym.edu.tw](mailto:alvin4016@ym.edu.tw))

Matlab進階程式設計與專題實作 (碩博)

授課教師：盧家鋒

- CV & Publications
- Members
- Research Interests
- Teaching Materials**
- Download Platforms
- Activities
- Relevant Links

- MRI (UG)
- MRM (UG)
- MRI Research (G)
- MATLAB programming (UG)
- MATLAB ML (G)**
- MATLAB GUI (G)
- Signal Processing (G)
- Computer Sci. (UG)
- Computer Arch. (UG)
- fMRI Analysis (G)
- rs-fMRI Analysis (G)
- fNIRS Basics (G)
- fNIRS Workshop (G)
- Human Dissection (UG)
- Neuroanatomy (UG)
- Image Processing (R)



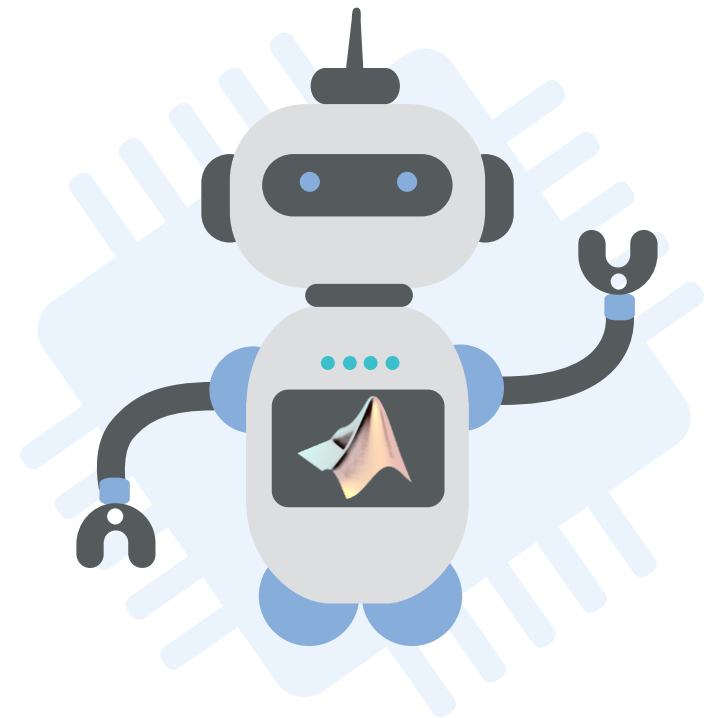
# Contents in this Week

## 01 Single-Layer Neural Network

Basic Concepts and Supervised Learning

## 02 Multi-Layer Neural Network

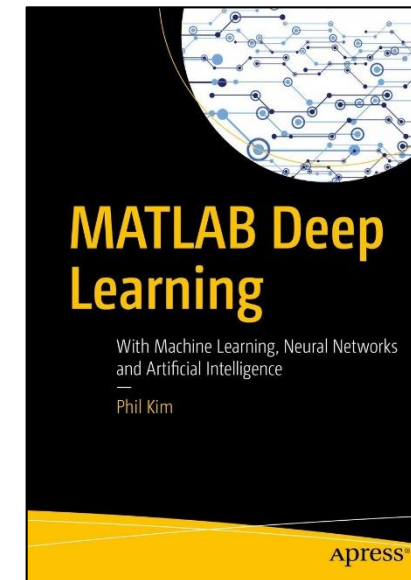
Back-propagation Algorithm, Momentum,  
Cross Entropy, Regularization

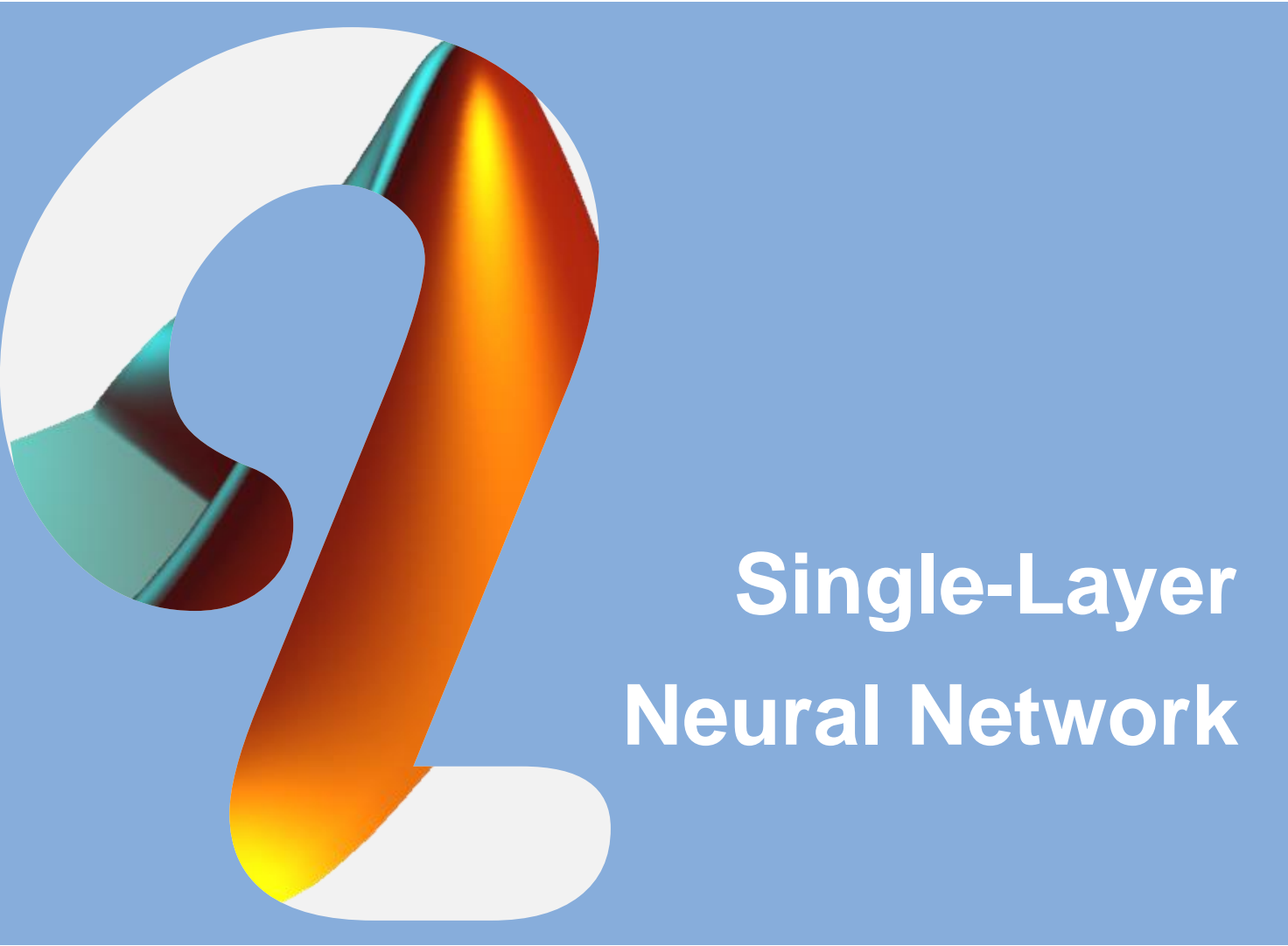


# References

## [Textbook 5]

- **MATLAB Deep Learning With Machine Learning, Neural Networks and Artificial Intelligence, 1st edition, 2017**  
Phil Kim
- **Neural Network (Ch.2)**
- **Training of Multi-Layer Neural Network (Ch.3)**





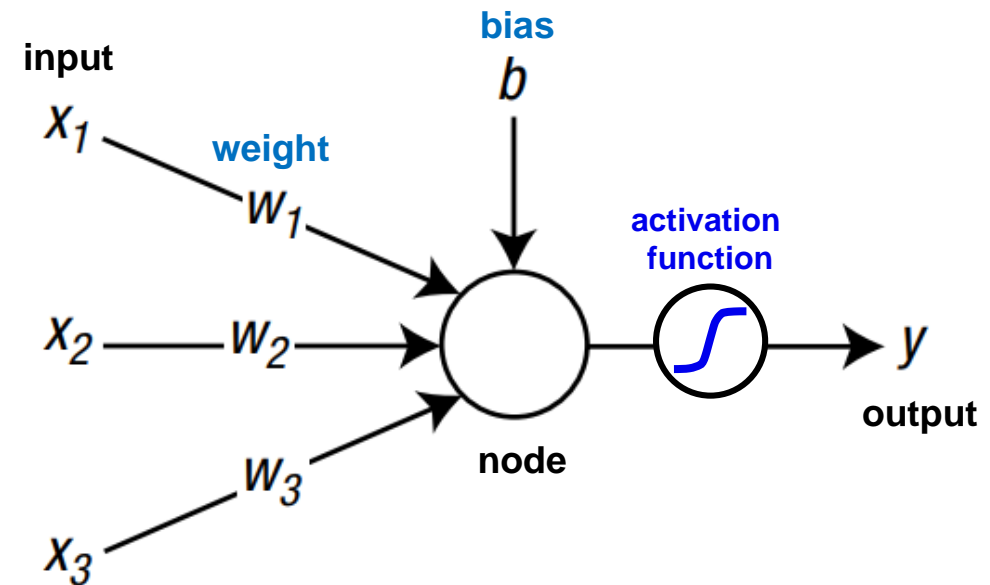
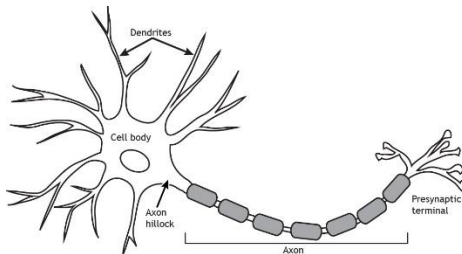
# Single-Layer Neural Network

Basic Concepts and  
Supervised Learning

# (Artificial) Neural Network

- The neural network imitates the mechanism of the brain. As the brain is composed of connections of numerous neurons
- The information of the neural net is stored in the form of **weights** and **bias**.

Brain	Neural Network
Neuron cell	Node
Connection of neurons	Connection weight
Action potential	Activation function

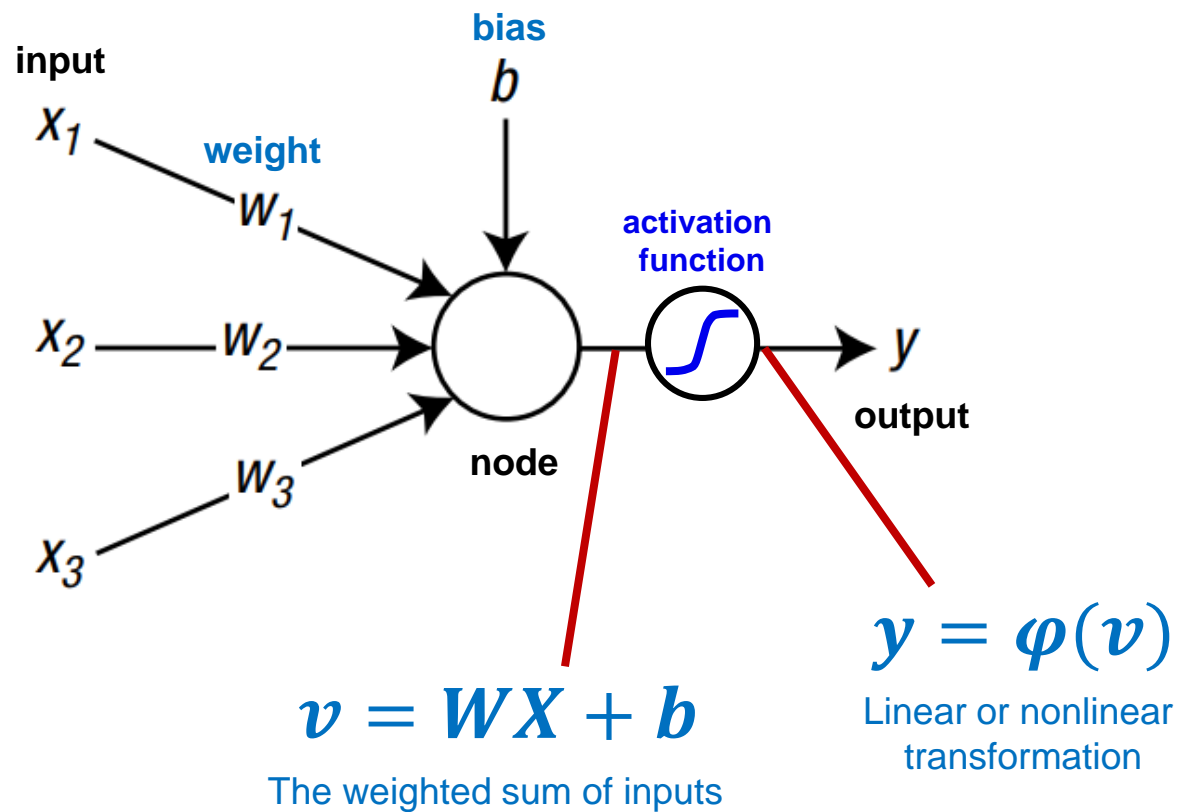


The arrows of the figure denote signal flow.

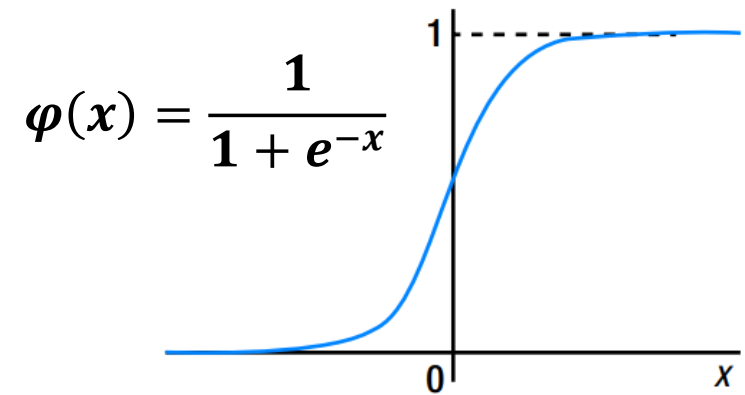
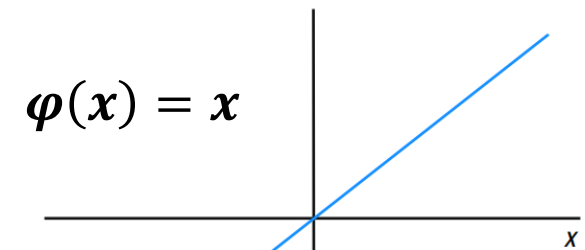


# Neural Network

Mathematical representations:

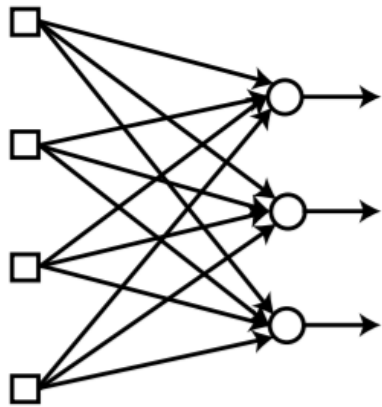


linear function

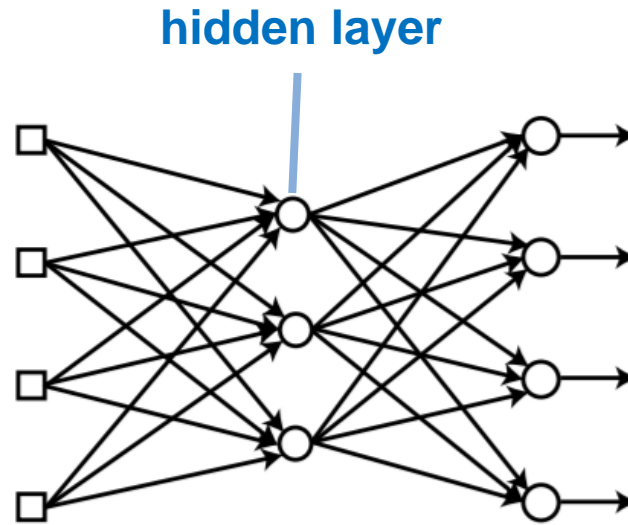


sigmoid function

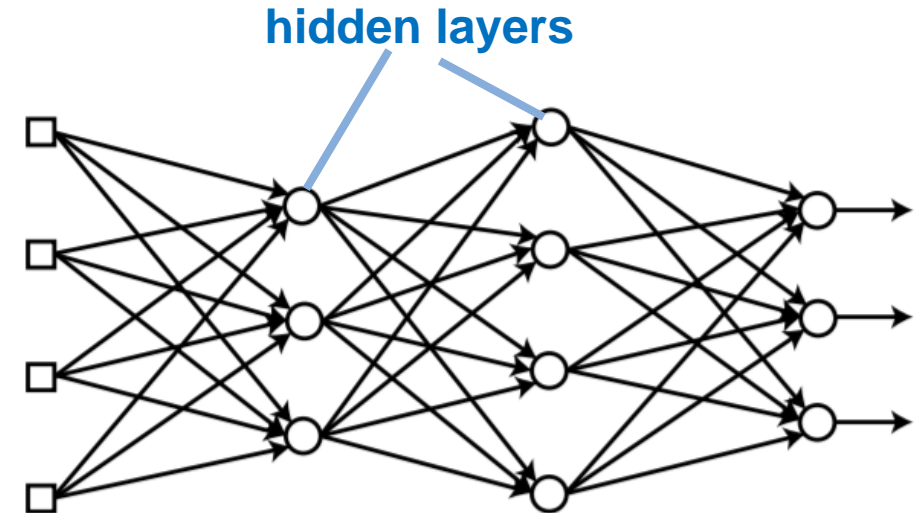
# Layers of Neural Network



Single-layer Neural Network



(Shallow) Multi-layer Neural Network



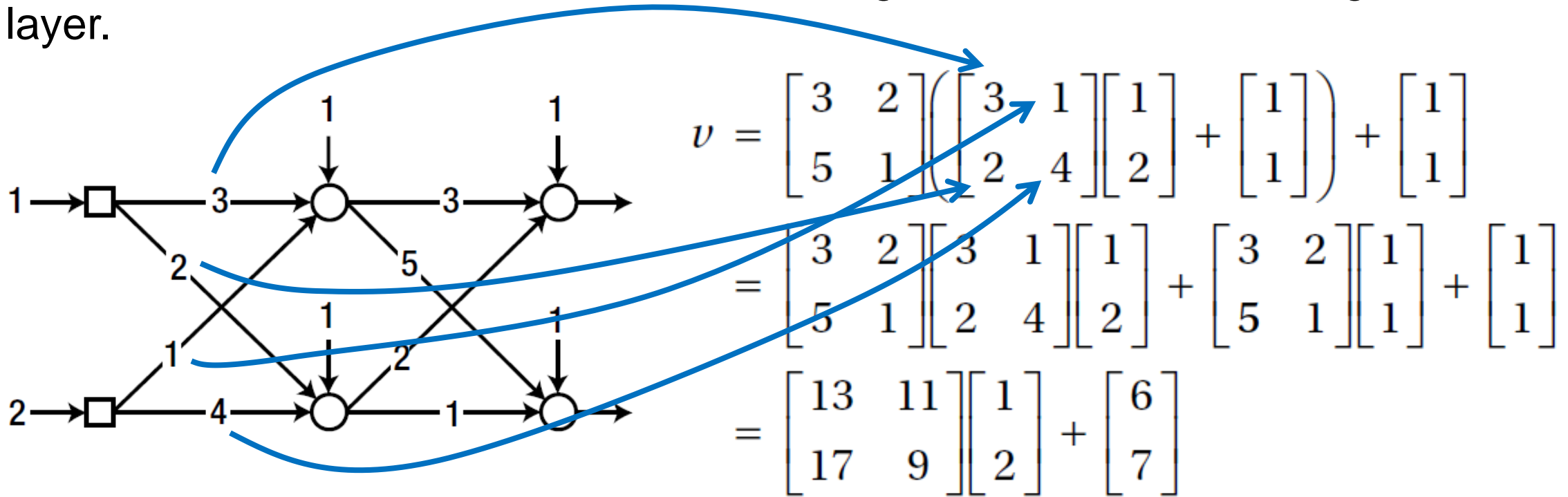
Deep Neural Network

Single-Layer Neural Network		Input Layer - Output Layer
Multi-Layer Neural Network	Shallow Neural Network	Input Layer - Hidden Layer - Output Layer
	Deep Neural Network	Input Layer - Hidden Layers - Output Layers



# Why nonlinear activation function?

- The use of a linear function for the nodes negates the effect of adding a layer.

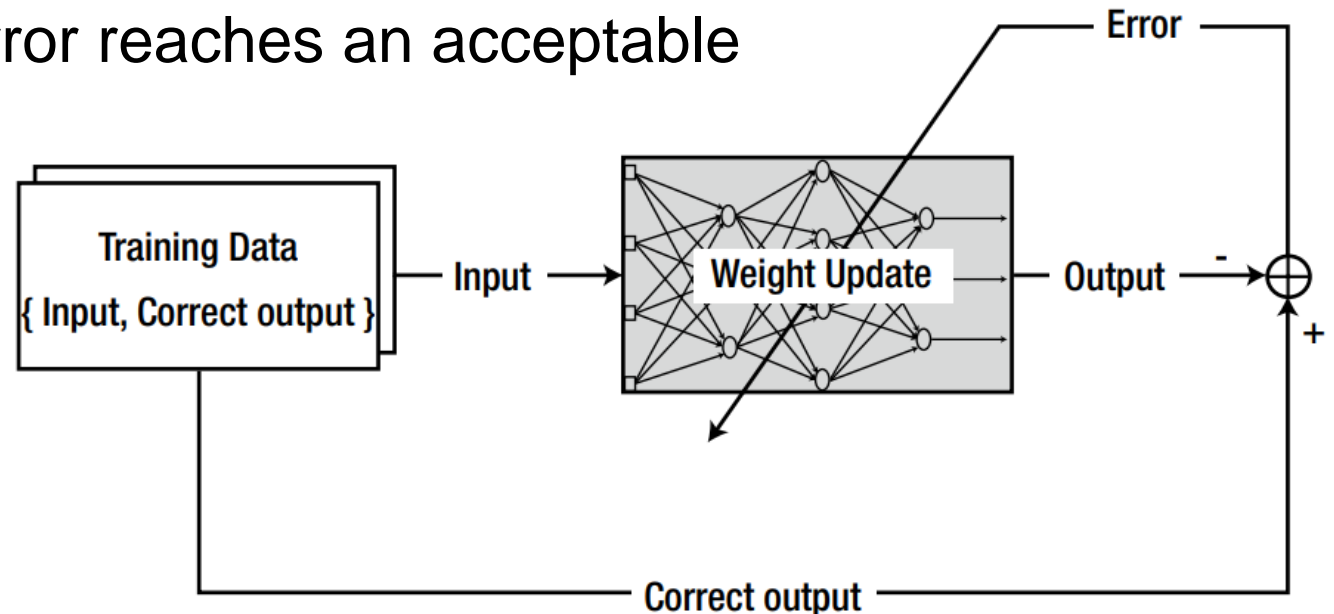


- In this case, the multi-layer model is mathematically identical to a single-layer neural network.

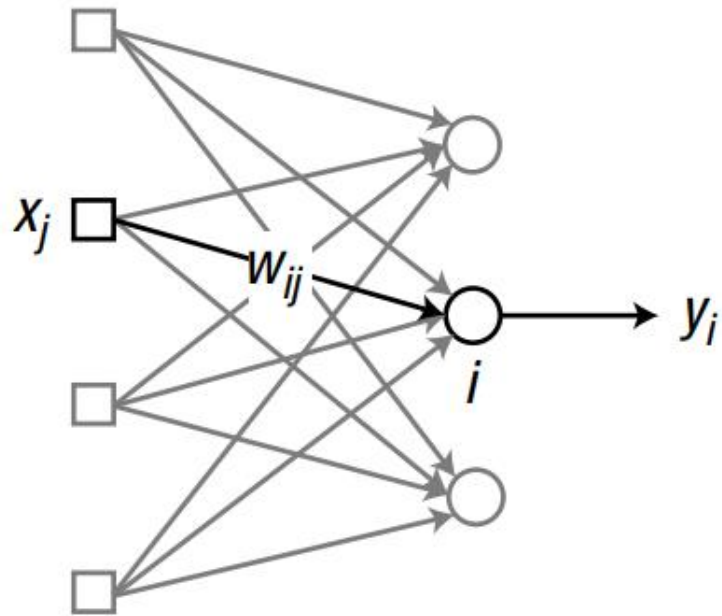
# Supervised Learning of a Neural Network

1. **Initialize** the weights with adequate values.
2. Enter the training data { input, correct output } into the neural network, and **calculate the error from the correct output**.
3. **Adjust the weights to reduce the error**.
4. Repeat Steps 2-3 for all training data
5. Repeat Steps 2-4 until the error reaches an acceptable tolerance level.

(All training data goes through Steps 2-4 once, is called an **epoch**.)



# Errors and Loss Function



$$e_i = d_i - y_i$$

$d_i$  is the **correct output** of the output node  $i$ .  
(ground truth)

- Let us define the loss function for output node  $y_i$

$$L_i = \frac{1}{2} (d_i - y_i)^2,$$

$$y_i = \varphi(v_i), \quad v_i = \sum_{j=1}^m w_{ij} x_j$$

where  $m$  is the numbers of input nodes

# Steepest Gradient Decent

- We minimize the loss function  $L_i$  w.r.t  $w_{ij}$

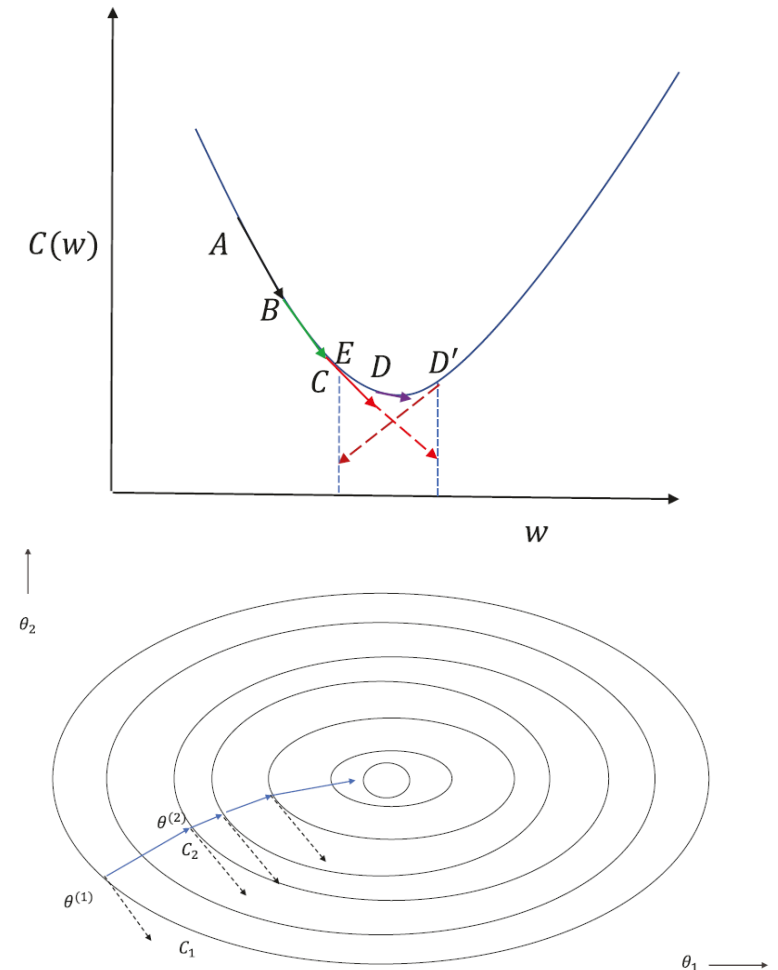
$$\frac{\partial L_i}{\partial w_{ij}} = e_i(-1) \frac{\partial \phi}{\partial v_i} \frac{\partial v_i}{\partial w_{ij}} = -e_i \phi' x_j$$

- The **steepest gradient decent** method

$$\begin{aligned} w_{ij}^{(k+1)} &= w_{ij}^{(k)} - \alpha \frac{\partial L}{\partial w_{ij}} \\ &= w_{ij}^{(k)} + \alpha \phi' e_i x_j \end{aligned}$$

- Or we may express the above equation as

$$w_{ij} \leftarrow w_{ij} + \alpha \phi' e_i x_j$$



# Generalized Delta Rule

- For an arbitrary activation function, the delta rule is expressed as

$$w_{ij} \leftarrow w_{ij} + \alpha \delta_i x_j$$

$$\delta_i = \varphi'(v_i) e_i$$

- The weight is adjusted in proportion to the input value,  $x_j$  and the output error,  $e_i$ .

$w_{ij}$  = The weight between the output node  $i$  and input node  $j$

$e_i$  = The error of the output node  $i$

$v_i$  = The weighted sum of the output node  $i$

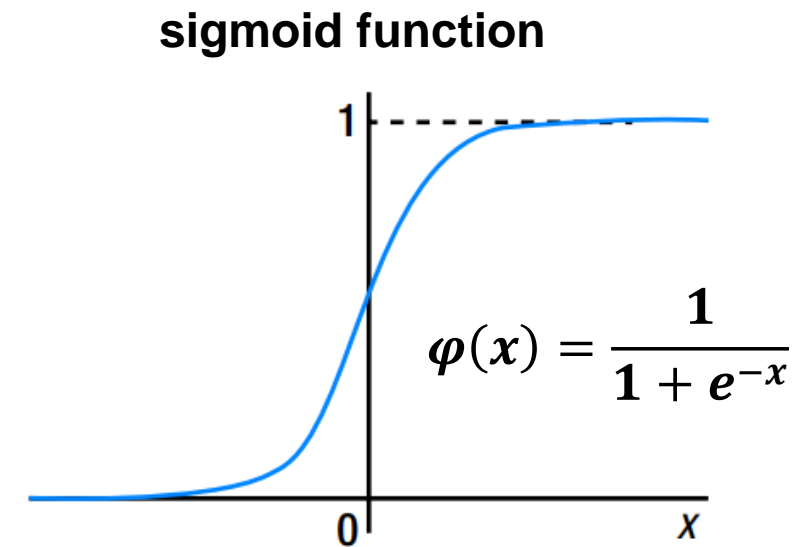
$\varphi'$  = The derivative of the activation function  $\varphi$  of the output node  $i$

$\alpha$  = Learning rate (  $0 < \alpha \leq 1$  )

# Derivative of sigmoid function

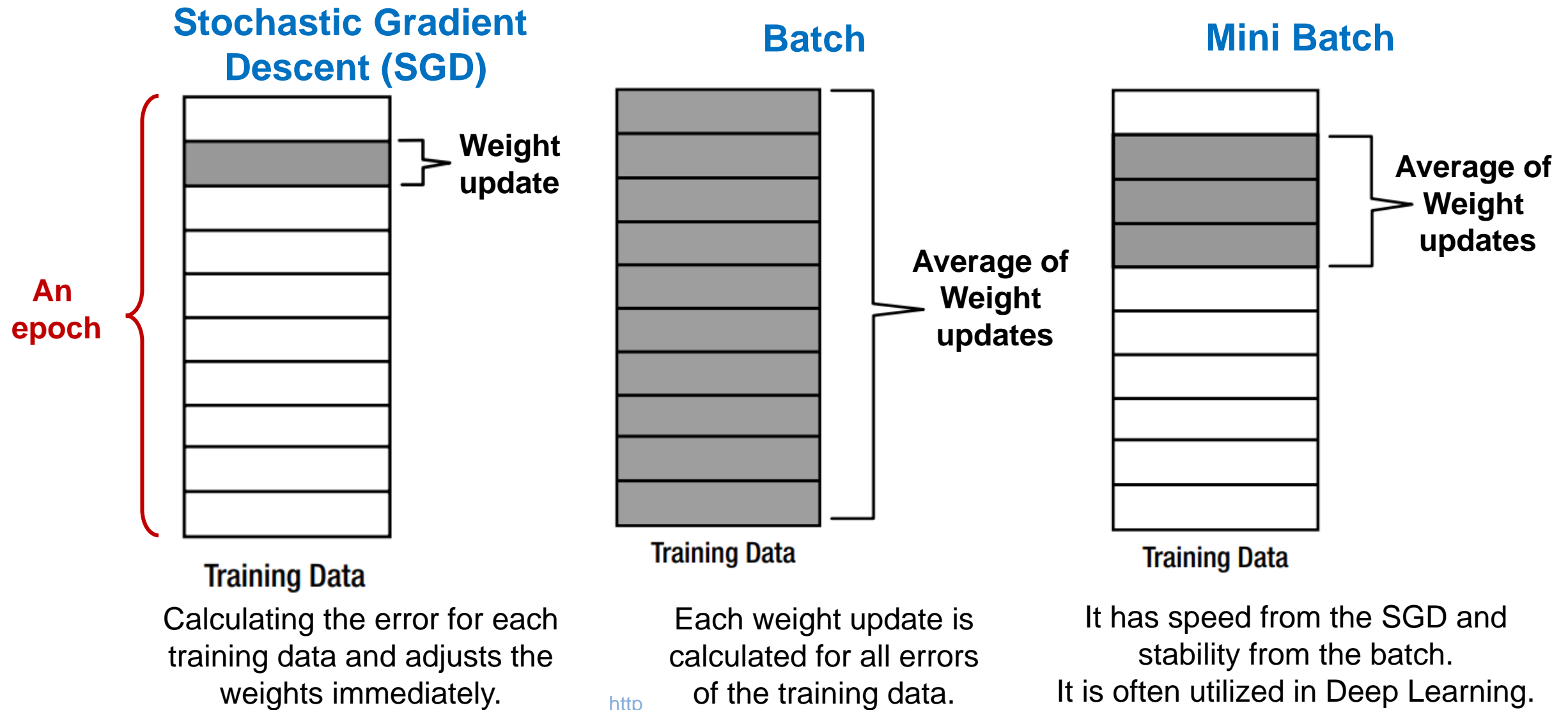
$$\begin{aligned}\varphi'(x) &= \frac{d(1 + e^{-x})^{-1}}{dx} = -(1 + e^{-x})^{-2}(-e^{-x}) \\ &= \frac{1}{1 + e^{-x}} \left( 1 - \frac{1}{1 + e^{-x}} \right) \\ &= \varphi(x)(1 - \varphi(x))\end{aligned}$$

$$w_{ij} \leftarrow w_{ij} + \alpha \varphi(x)(1 - \varphi(x)) e_i x_j$$





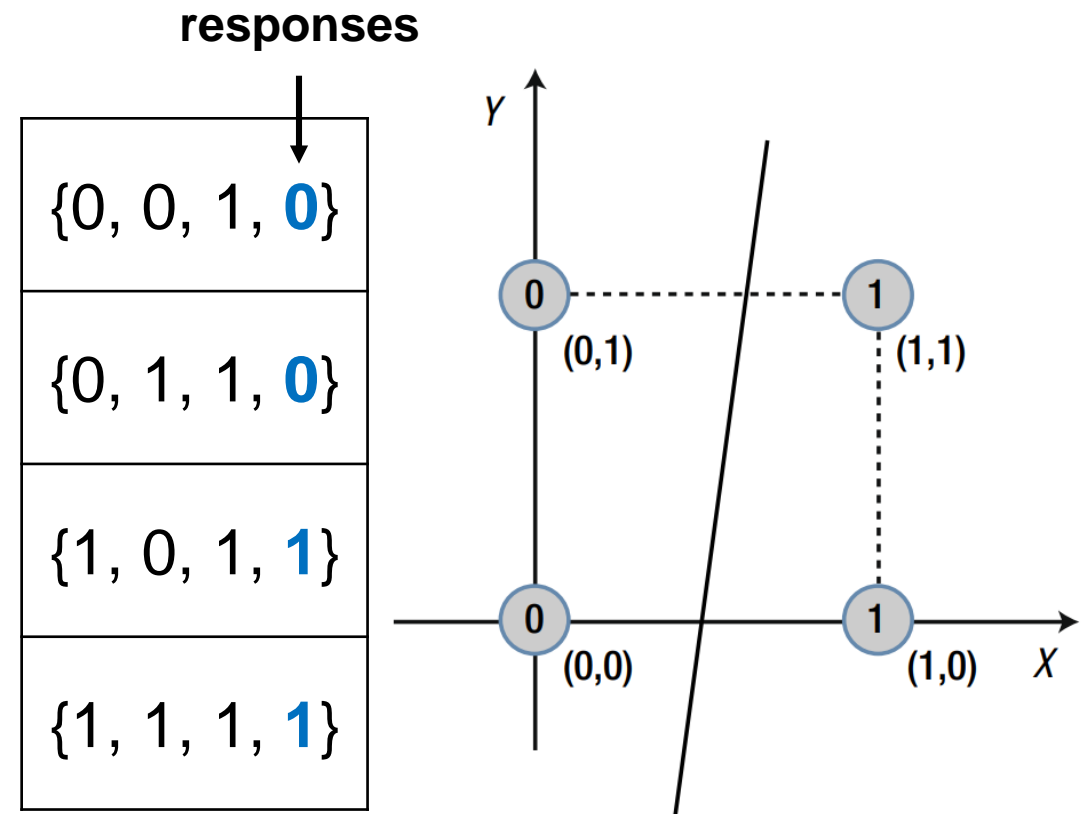
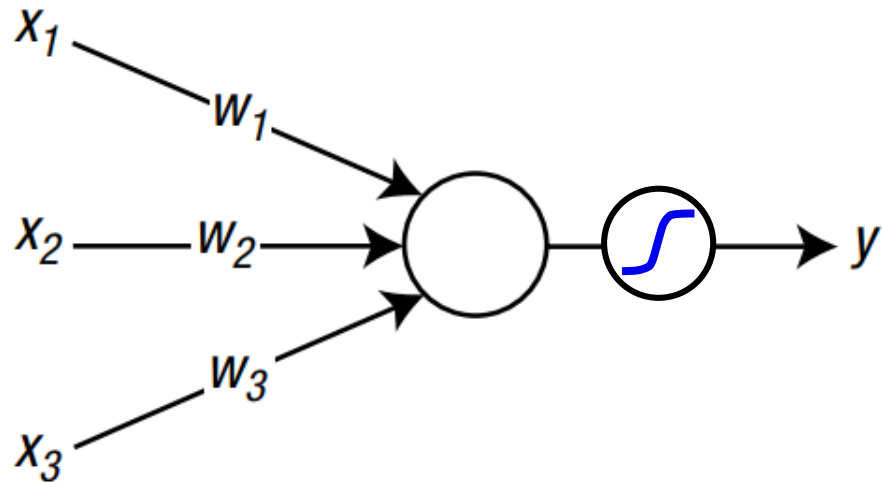
# Schemes for Updating weights



# Example 1: Linearly Separable

- **MLmaterials\_L10\Single-layer\**

- **TestDeltaSGD.m**
- **DeltaSGD.m**
- **Sigmoid.m**



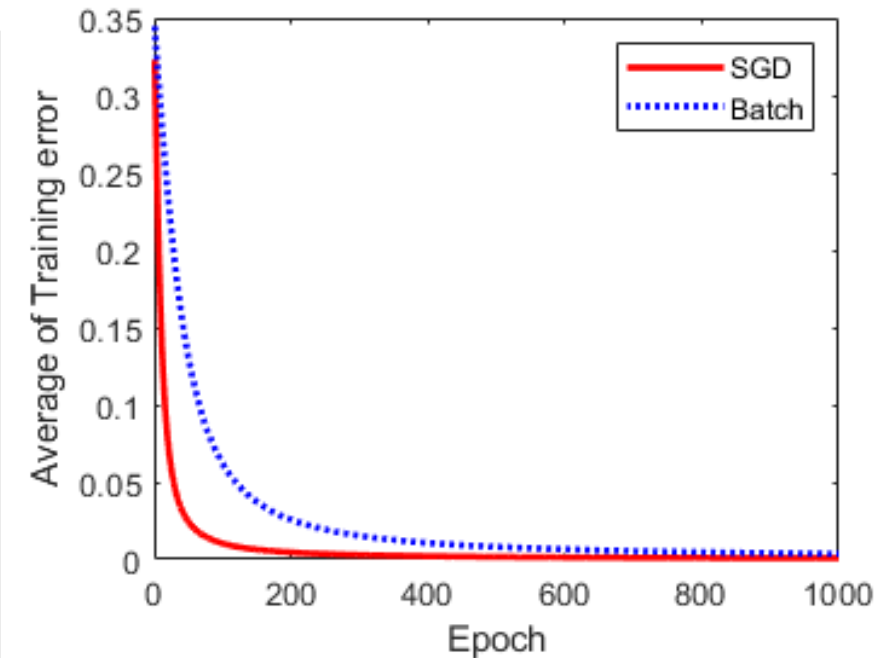
# SGD vs. Batch

- TestDeltaSGD.m
- DeltaSGD.m

```
9 -      v = W*x;  
10 -      y = Sigmoid(v);  
11  
12 -      e      = d - y;  
13 -      delta = y*(1-y)*e;  
14  
15 -      dW = alpha*delta*x;  
16  
17 -      W(1) = W(1) + dW(1);  
18 -      W(2) = W(2) + dW(2);  
19 -      W(3) = W(3) + dW(3);  
20 -      end  
21 -      end
```

- TestDeltaBatch.m
- DeltaBatch.m

```
11 -      v = W*x;  
12 -      y = Sigmoid(v);  
13  
14 -      e      = d - y;  
15 -      delta = y*(1-y)*e;  
16  
17 -      dW = alpha*delta*x;  
18  
19 -      dWsum = dWsum + dW;  
20 -      end  
21 -      dWavg = dWsum / N;  
22  
23 -      W(1) = W(1) + dWavg(1);  
24 -      W(2) = W(2) + dWavg(2);  
25 -      W(3) = W(3) + dWavg(3);  
26 -      end
```

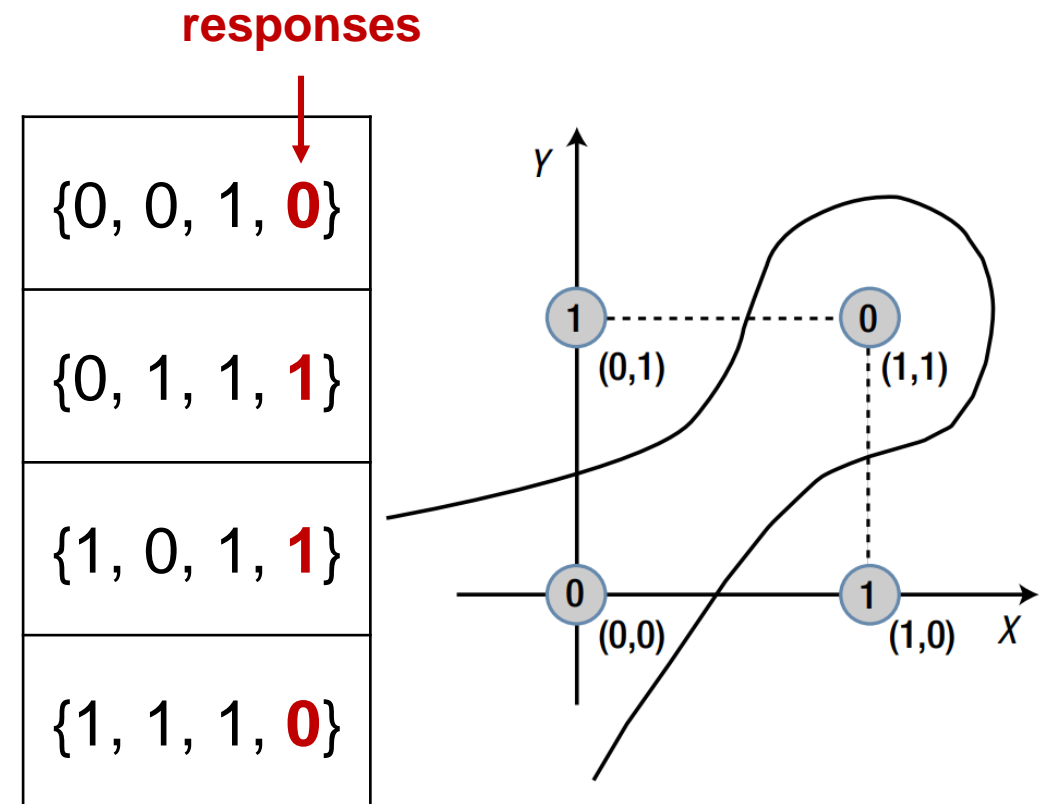
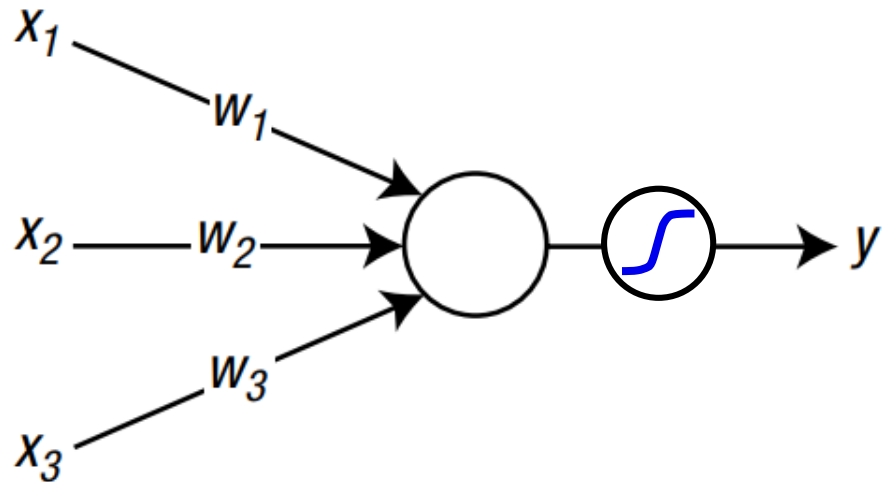


**Average of  
Weight  
updates**

# Example 2: Linearly Inseparable

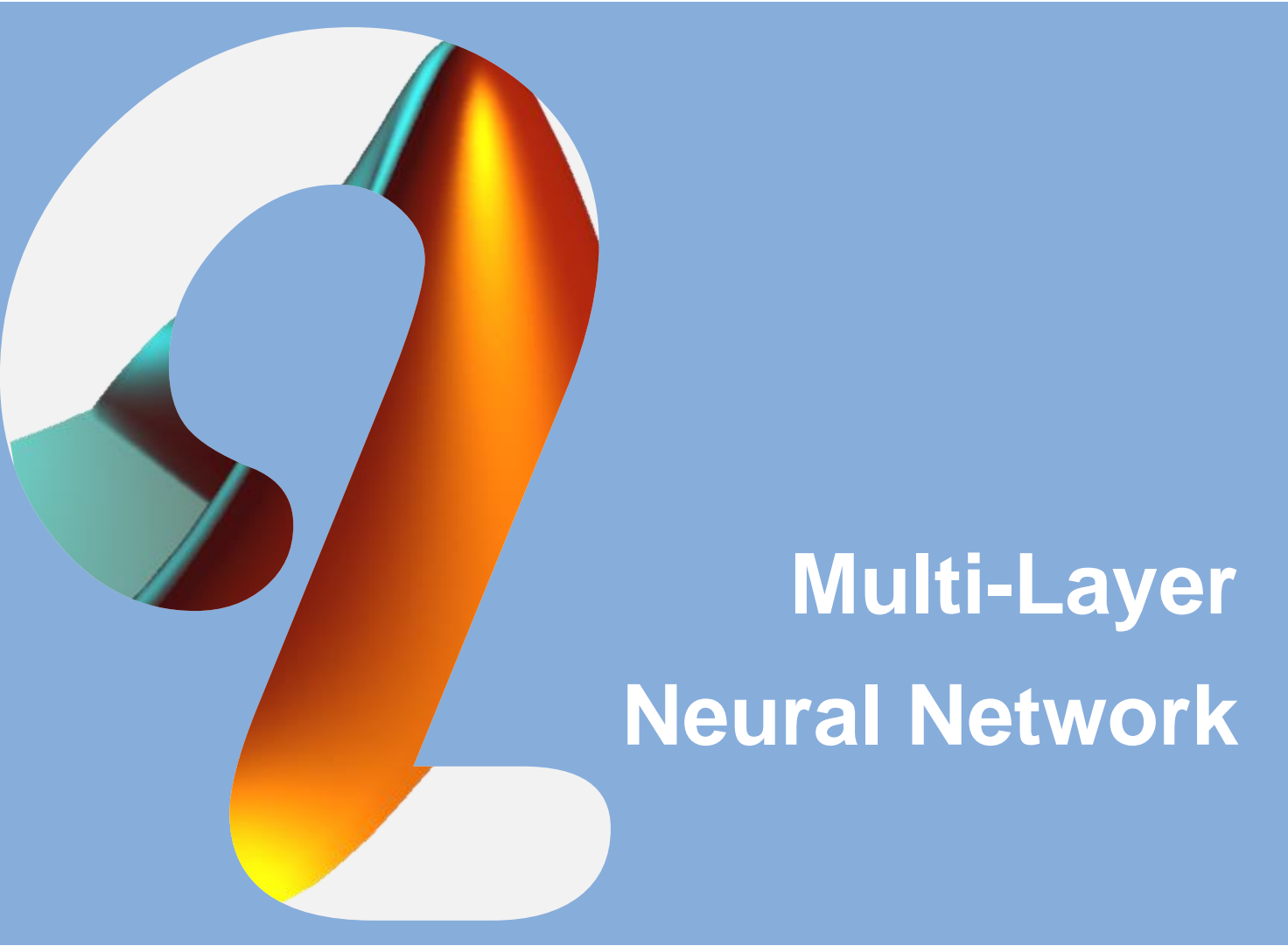
- MLmaterials\_L10\Single-layer\

- TestDeltaSGD.m
- DeltaSGD.m
- Sigmoid.m



# Short Summary

- The single-layer neural network can only solve linearly separable problems. This is because the single-layer neural network is a model that linearly divides the input data space.
- In order to overcome this limitation of the single-layer neural network, **we need more layers** in the network.



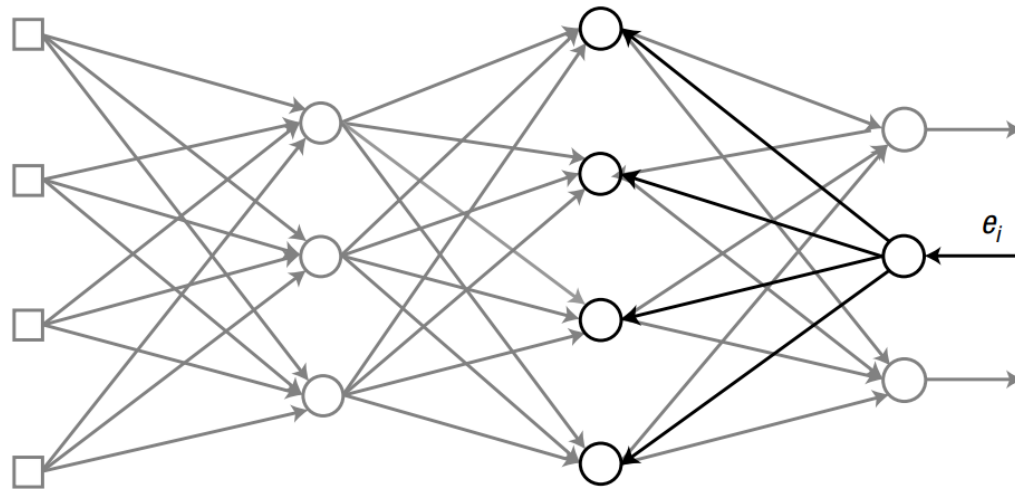
# Multi-Layer Neural Network

Back-propagation Algorithm,  
Momentum, Cross Entropy,  
Regularization



# Back-Propagation Algorithm

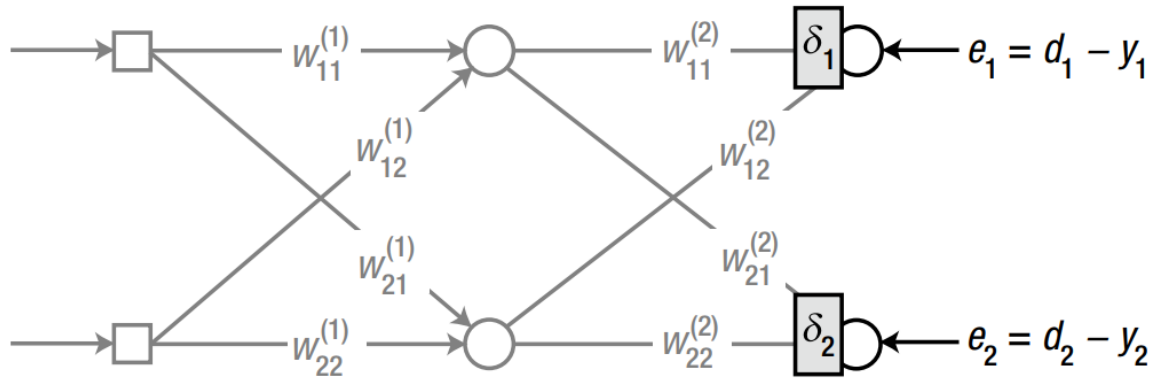
- The previously introduced **delta rule** is ineffective for training of the multi-layer neural network because the error is not defined in the hidden layers.
- **Back-propagation algorithm** provided a systematic method to determine the error of the hidden nodes.
- Once the hidden layer errors are determined, the delta rule is applied to adjust the weights.



the output error starts from the output layer and **moves backward** until it reaches the right next hidden layer to the input layer.

# Back-Propagation Algorithm

- The first thing to calculate is delta,  $\delta$ , of each node :



$$e_1 = d_1 - y_1$$

$$\delta_1 = \varphi'(v_1)e_1$$

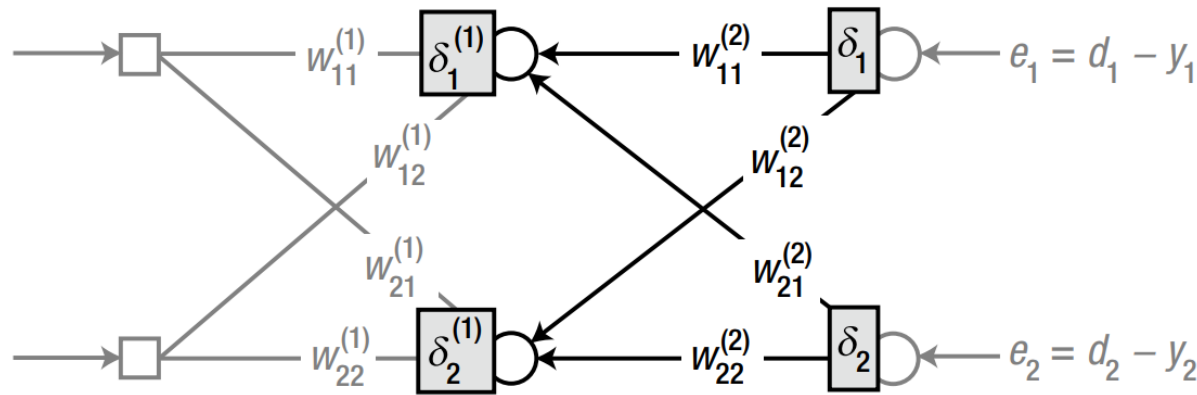
$$e_2 = d_2 - y_2$$

$$\delta_2 = \varphi'(v_2)e_2$$

$\varphi'$  is the derivative of the activation function of the output node.  
 $y_i$  is the output from the output node.  
 $d_i$  is the correct output from the training data.  
 $v_i$  is the weighted sum of the corresponding node.

# Back-Propagation Algorithm

- Since we have  $\delta_1$  and  $\delta_2$ , let's proceed leftward to the hidden nodes and calculate the delta :



$$e_1^{(1)} = w_{11}^{(2)} \delta_1 + w_{21}^{(2)} \delta_2$$

$$\begin{aligned} \delta_1^{(1)} &= \varphi' \left( v_1^{(1)} \right) e_1^{(1)} \\ &= \varphi \left( v_1^{(1)} \right) (1 - \varphi \left( v_1^{(1)} \right)) e_1^{(1)} \end{aligned}$$

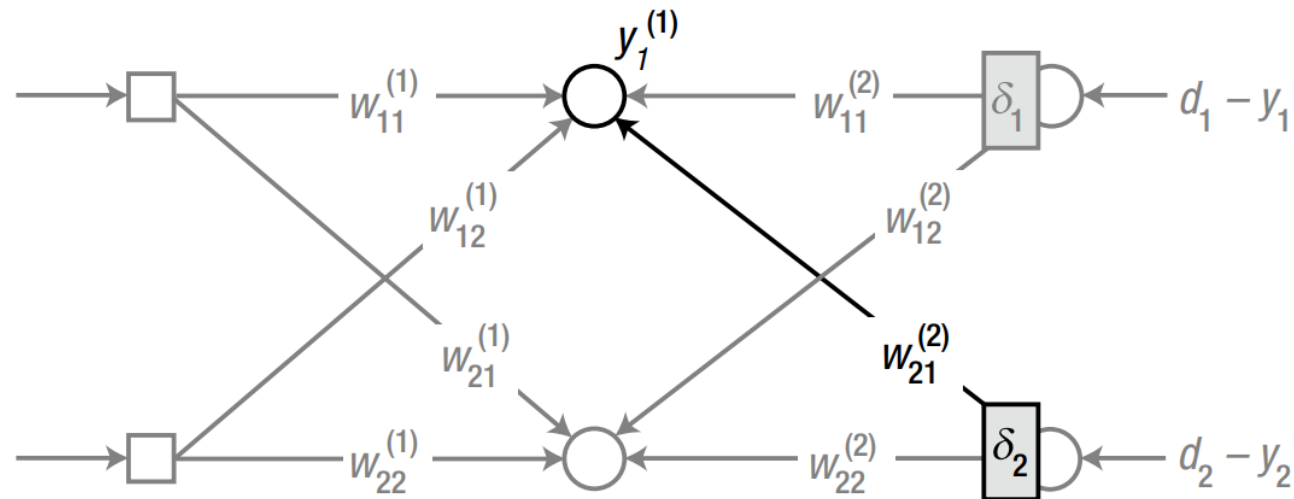
$$e_2^{(1)} = w_{12}^{(2)} \delta_1 + w_{22}^{(2)} \delta_2$$

$$\begin{aligned} \delta_2^{(1)} &= \varphi' \left( v_2^{(1)} \right) e_2^{(1)} \\ &= \varphi \left( v_2^{(1)} \right) (1 - \varphi \left( v_2^{(1)} \right)) e_2^{(1)} \end{aligned}$$

$v_1^{(1)}$  and  $v_2^{(1)}$  are the weighted sums of the forward signals at the respective nodes.

# Update of Weights

- Consider the weight  $w_{21}^{(2)}$  for example.

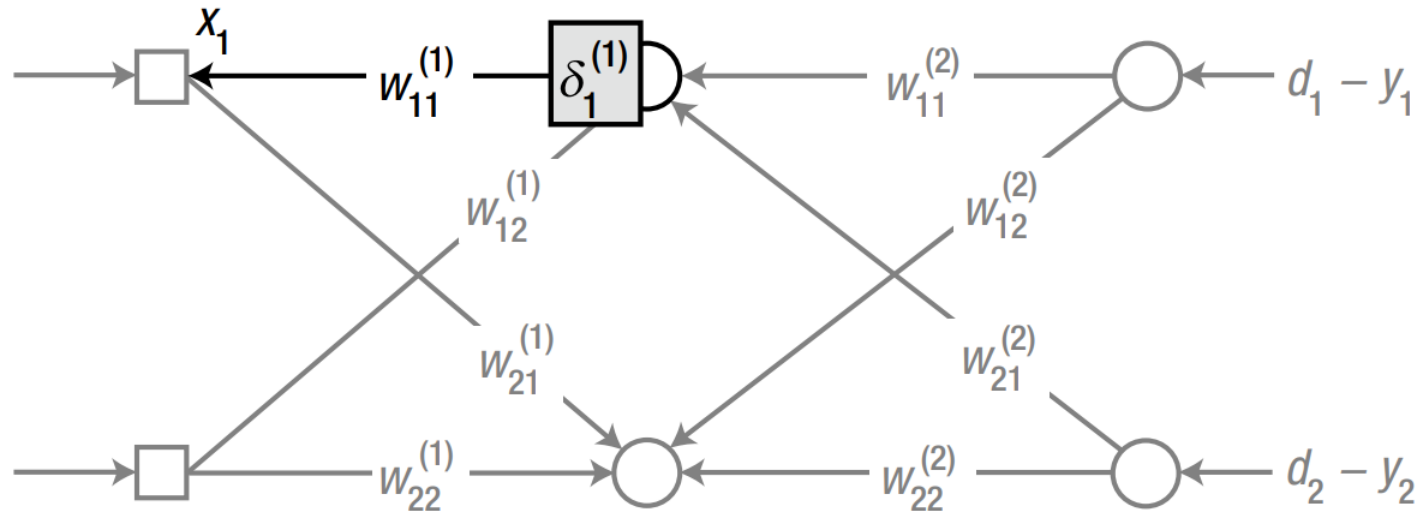


$$w_{21}^{(2)} \leftarrow w_{21}^{(2)} + \alpha \delta_2 y_1^{(1)}$$

$y_1^{(1)}$  is the output of the first hidden node.

# Update of Weights

- The weight  $w_{11}^{(1)}$  of figure is adjusted as :

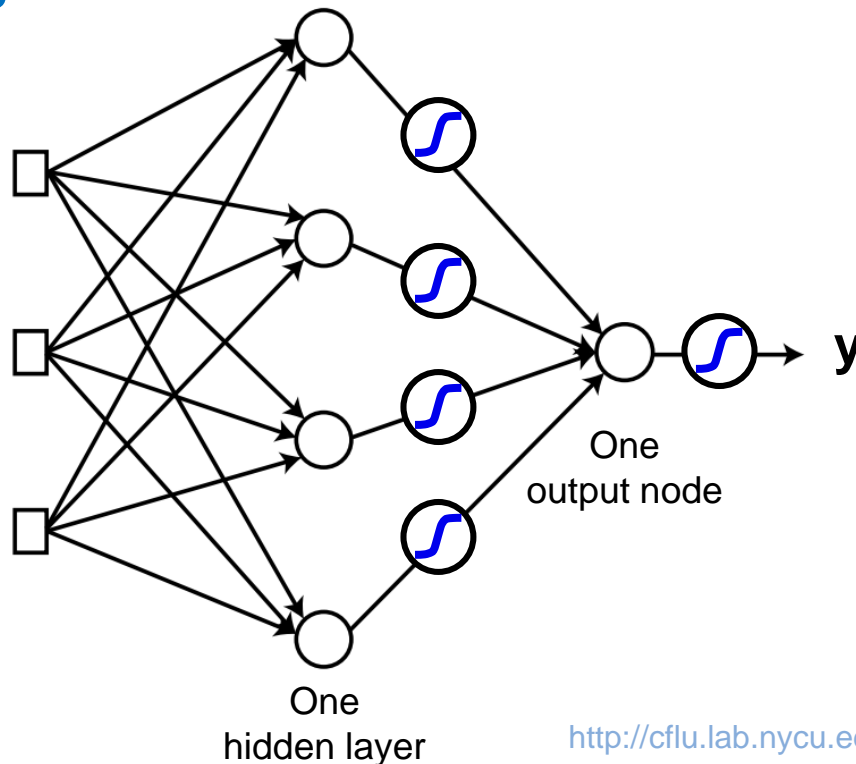


$$w_{11}^{(1)} \leftarrow w_{11}^{(1)} + \alpha \delta_1^{(1)} x_1$$

$x_1$  is the output of the first input node.

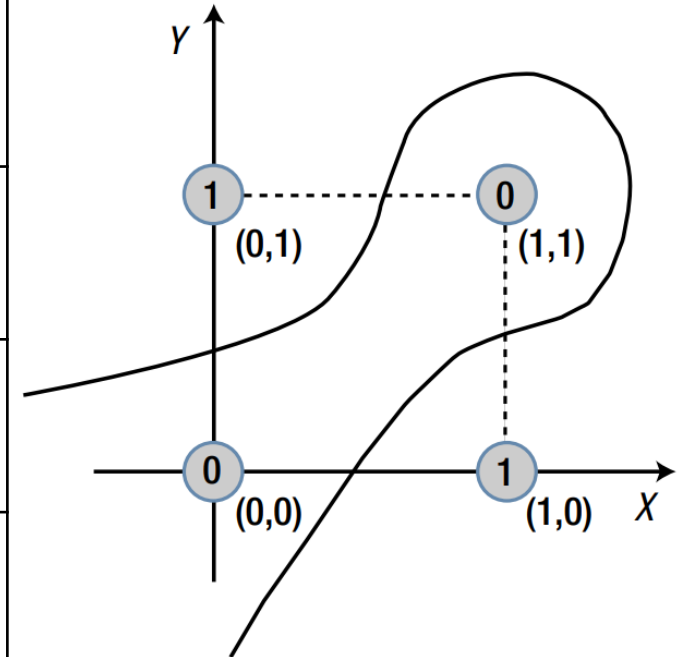
# Example: Back-Propagation

- **MLmaterials\_L10\Multi-layer\**
  - **TestBackpropXOR.m**
  - **BackpropXOR.m**
  - **Sigmoid.m**



responses

$\{0, 0, 1, \mathbf{0}\}$
$\{0, 1, 1, \mathbf{1}\}$
$\{1, 0, 1, \mathbf{1}\}$
$\{1, 1, 1, \mathbf{0}\}$





# Momentum

- The momentum,  $m$ , is a term that is added to the delta rule for adjusting the weight.
- It acts similarly to physical momentum, which impedes the reaction of the body to the external forces.

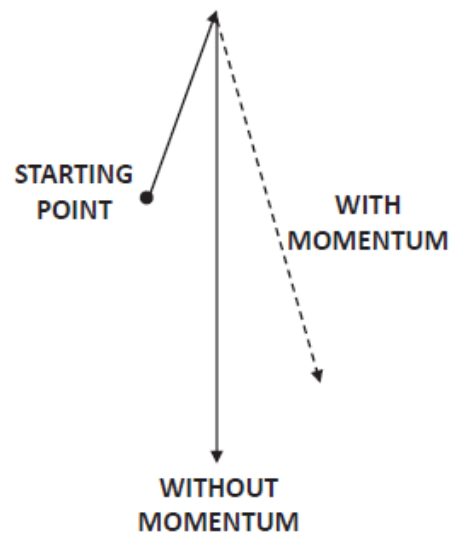
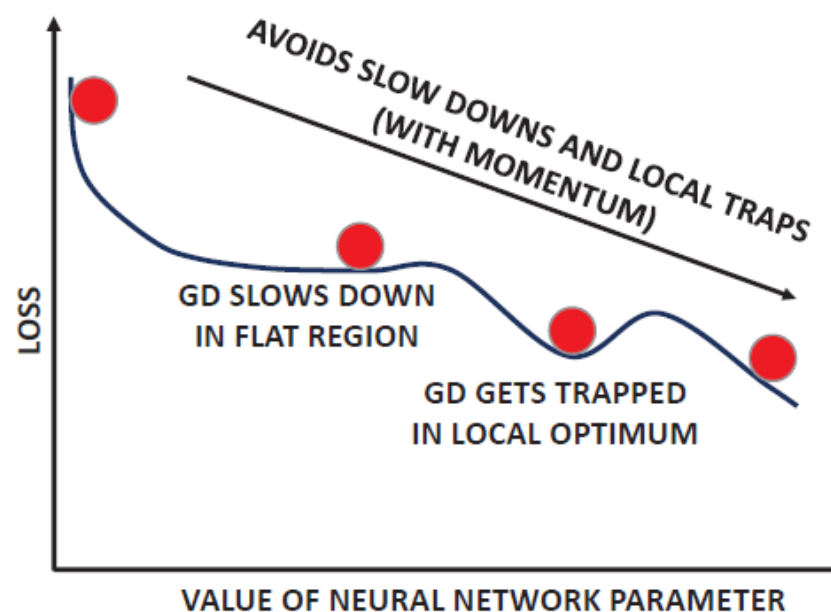
$$\begin{array}{l} \Delta w_{ij} = \alpha \delta_i x_j \\ w_{ij} \leftarrow w_{ij} + \Delta w_{ij} \end{array} \quad \Rightarrow \quad \begin{array}{l} \Delta w_{ij} = \alpha \delta_i x_j \\ m = \Delta w_{ij} + \beta m^- \\ w_{ij} \leftarrow w_{ij} + m \\ m^- = m \end{array}$$

$m^-$  is the previous momentum

$\beta$  is a positive constant that is less than 1.

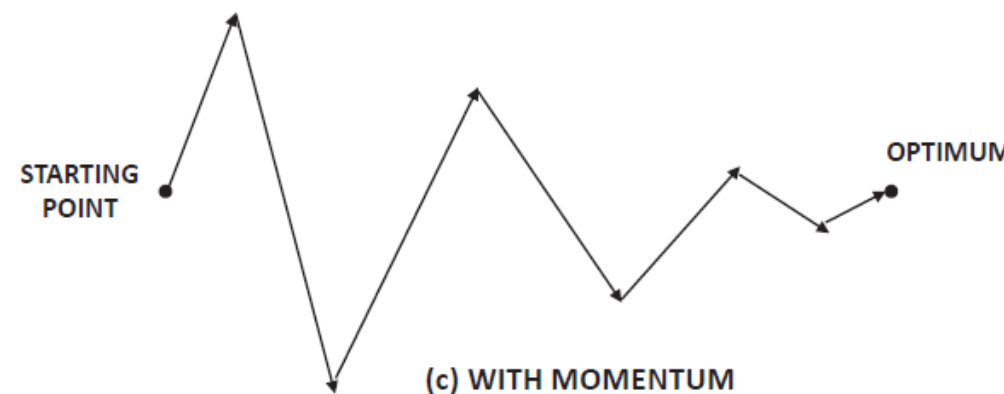
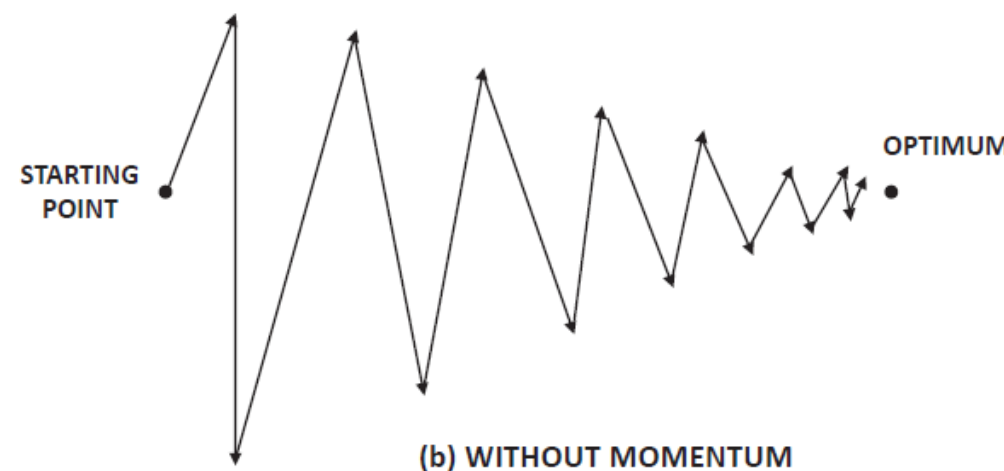
# Momentum improves the learning stability.

## Marble Rolling Down Hill



(a) RELATIVE DIRECTIONS

## Avoiding Zig-Zagging with Momentum

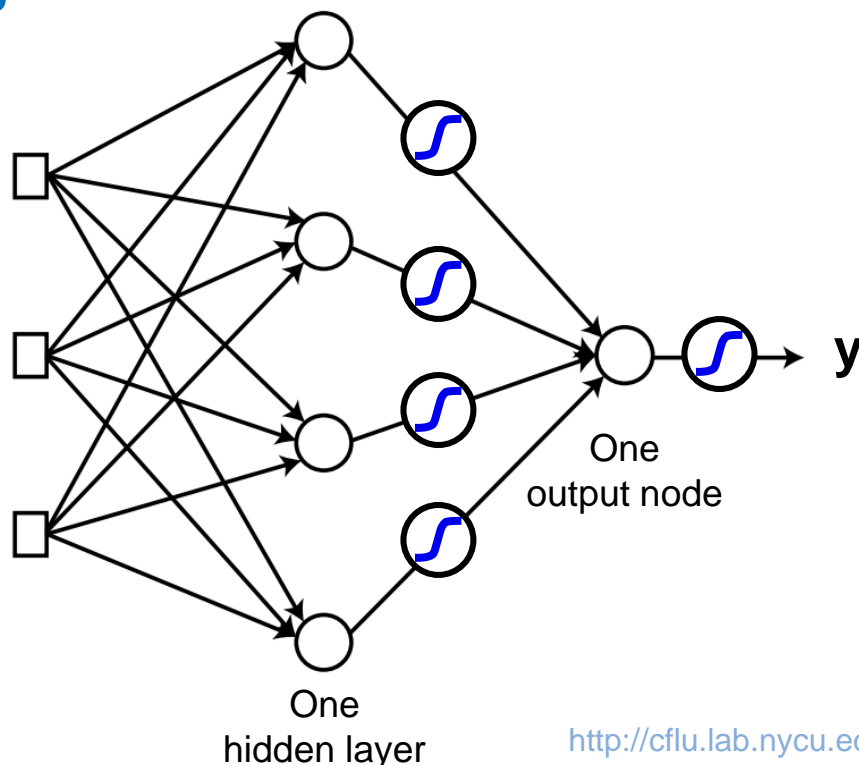


Charu C. Aggarwal, Neural Networks and Deep Learning, Springer, 2018

# Example: Back-Propagation with Momentum

- **MLmaterials\_L10\Multi-layer\**

- **TestBackpropMmt.m**
- **BackpropMmt.m**
- **Sigmoid.m**



responses

$\{0, 0, 1, \mathbf{0}\}$
$\{0, 1, 1, \mathbf{1}\}$
$\{1, 0, 1, \mathbf{1}\}$
$\{1, 1, 1, \mathbf{0}\}$

with momentum

$$\begin{bmatrix} 0.0038 \\ 0.9929 \\ 0.9919 \\ 0.0127 \end{bmatrix}$$

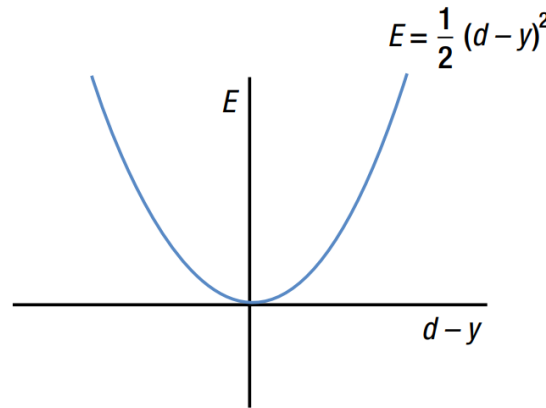
without momentum

$$\begin{bmatrix} 0.0060 \\ 0.9888 \\ 0.9891 \\ 0.0134 \end{bmatrix}$$

# Cost Function and Learning Rule

- There are two primary types of cost functions

$$L = \sum_{i=1}^M \frac{1}{2} (d_i - y_i)^2$$

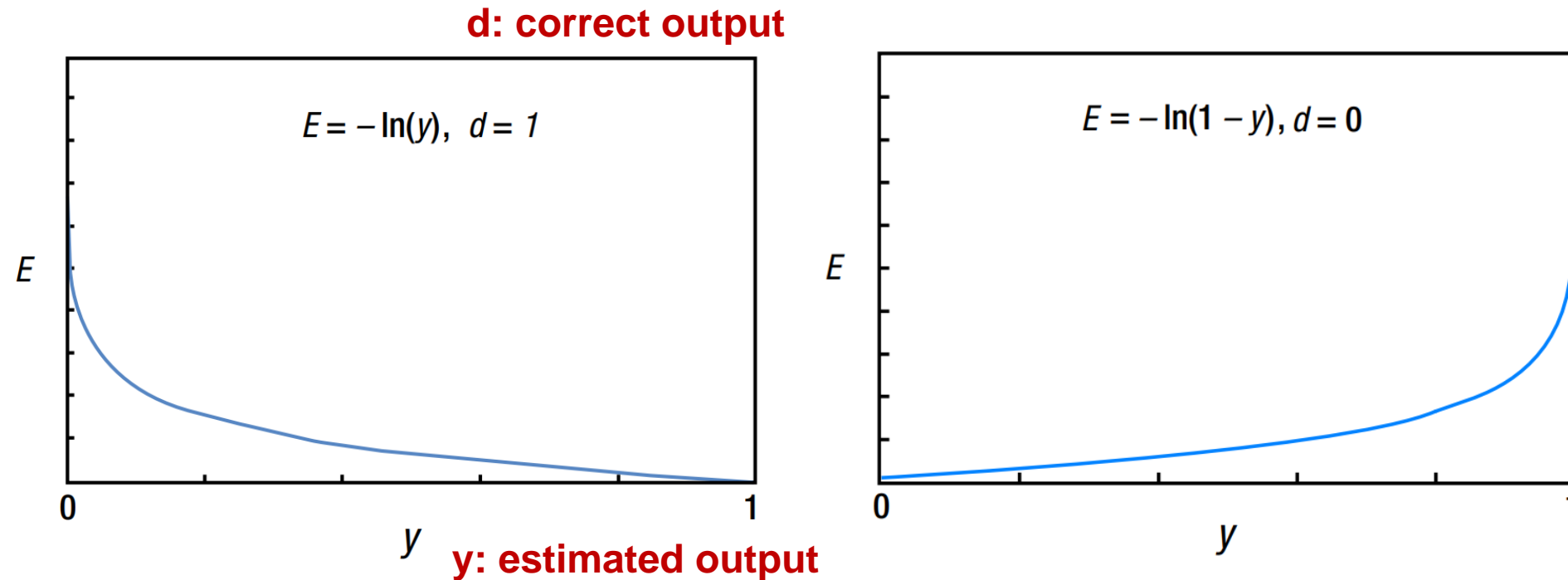


**Quadratic function**

$$L = \sum_{i=1}^M \{-d_i \ln(y_i) - (1 - d_i) \ln(1 - y_i)\}$$

**Cross entropy function**  
(tend to have better performance)

# Cross Entropy Function



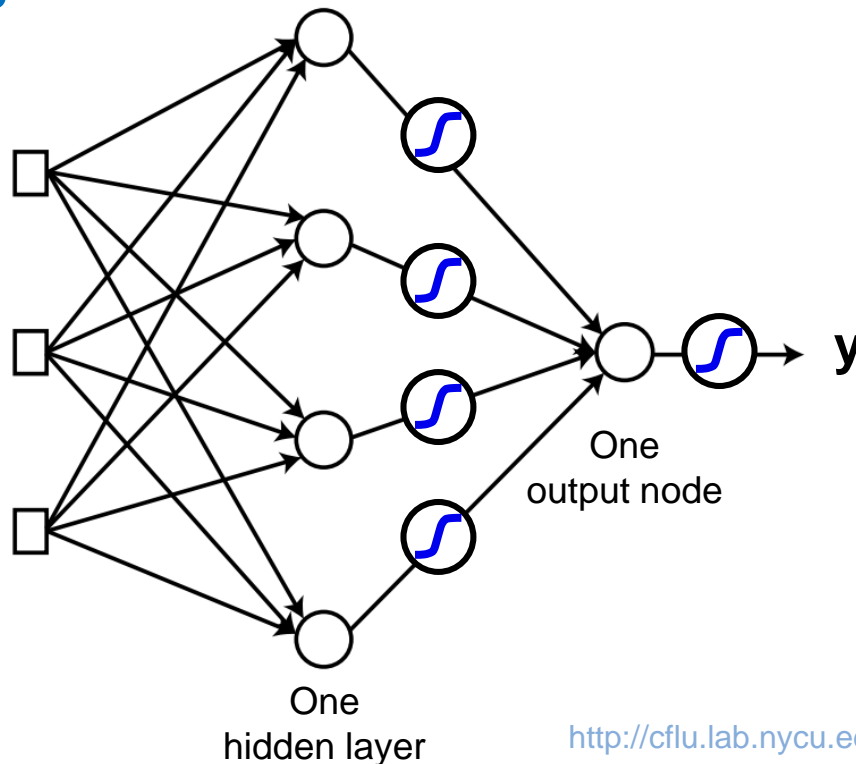
- This cost function is proportional to the error.
- The cross entropy function is much more sensitive to the error than quadratic function.

# Example:

## Back-Propagation using cross entropy

- **MLmaterials\_L10\Multi-layer\**

- **TestBackpropCE.m**
- **BackpropCE.m**
- **Sigmoid.m**



responses

$\{0, 0, 1, \mathbf{0}\}$
$\{0, 1, 1, \mathbf{1}\}$
$\{1, 0, 1, \mathbf{1}\}$
$\{1, 1, 1, \mathbf{0}\}$

Using cross entropy

$$\begin{bmatrix} 0.00003 \\ 0.9999 \\ 0.9998 \\ 0.00036 \end{bmatrix}$$

with momentum

$$\begin{bmatrix} 0.0038 \\ 0.9929 \\ 0.9919 \\ 0.0127 \end{bmatrix}$$



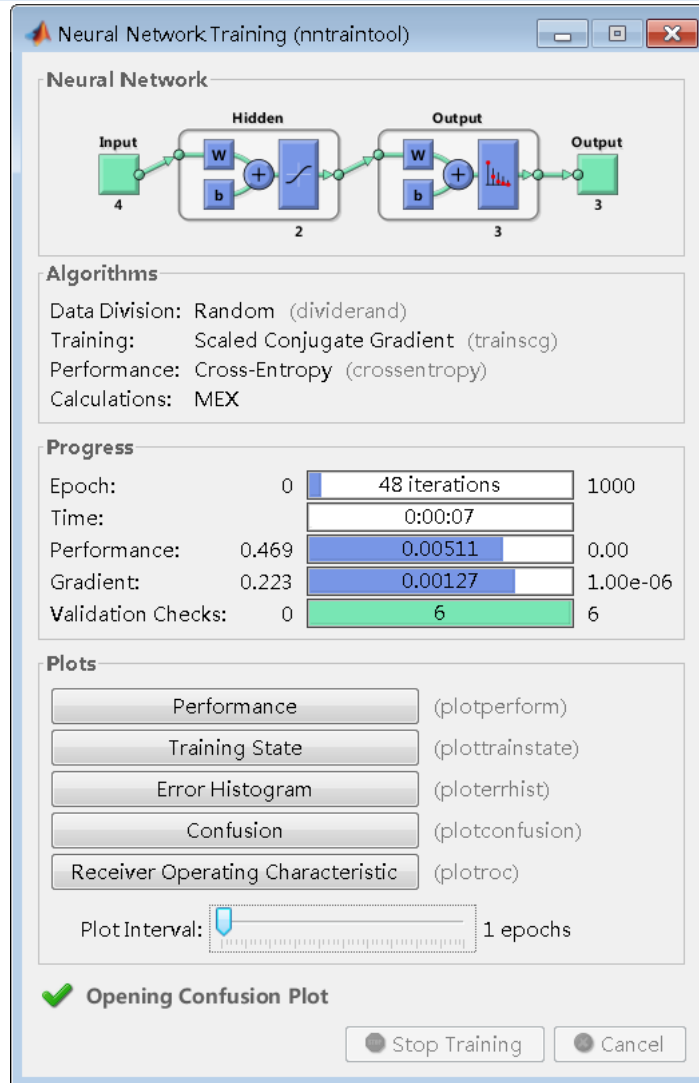
# Regularization

- One of the primary approaches used to overcome overfitting is making the model as simple as possible using regularization.
- In a mathematical sense, the essence of regularization is adding the sum of the weights to the cost function.

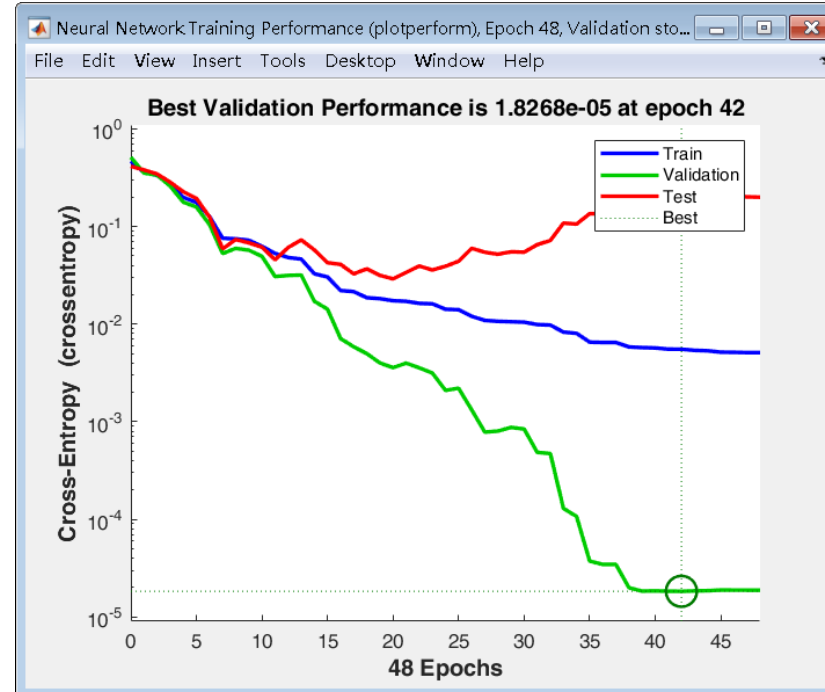
$$J = \frac{1}{2} \sum_{i=1}^M (d_i - y_i)^2 + \lambda \frac{1}{2} \|\mathbf{w}\|^2$$

$$J = \sum_{i=1}^M \{-d_i \ln(y_i) - (1 - d_i) \ln(1 - y_i)\} + \lambda \frac{1}{2} \|\mathbf{w}\|^2$$

# nnstart



## Iris dataset



**Training Confusion Matrix**

Output Class \ Target Class	1	2	3	
1	35 33.7%	0 0.0%	0 0.0%	100% 0.0%
2	0 0.0%	33 31.7%	0 0.0%	100% 0.0%
3	0 0.0%	1 1.0%	35 33.7%	97.2% 2.8%
	100% 0.0%	97.1% 2.9%	100% 0.0%	99.0% 1.0%

**Validation Confusion Matrix**

Output Class \ Target Class	1	2	3	
1	8 34.8%	0 0.0%	0 0.0%	100% 0.0%
2	0 0.0%	6 26.1%	0 0.0%	100% 0.0%
3	0 0.0%	0 0.0%	9 39.1%	100% 0.0%
	100% 0.0%	100% 0.0%	100% 0.0%	100% 0.0%

**Test Confusion Matrix**

Output Class \ Target Class	1	2	3	
1	7 30.4%	0 0.0%	0 0.0%	100% 0.0%
2	0 0.0%	10 43.5%	2 8.7%	83.3% 16.7%
3	0 0.0%	0 0.0%	4 17.4%	100% 0.0%
	100% 0.0%	100% 0.0%	66.7% 33.3%	91.3% 8.7%

**All Confusion Matrix**

Output Class \ Target Class	1	2	3	
1	50 33.3%	0 0.0%	0 0.0%	100% 0.0%
2	0 0.0%	49 32.7%	2 1.3%	96.1% 3.9%
3	0 0.0%	1 0.7%	48 32.0%	98.0% 2.0%
	100% 0.0%	98.0% 2.0%	96.0% 4.0%	98.0% 2.0%



THE END

**Contact:**

盧家鋒 alvin4016@nycu.edu.tw