



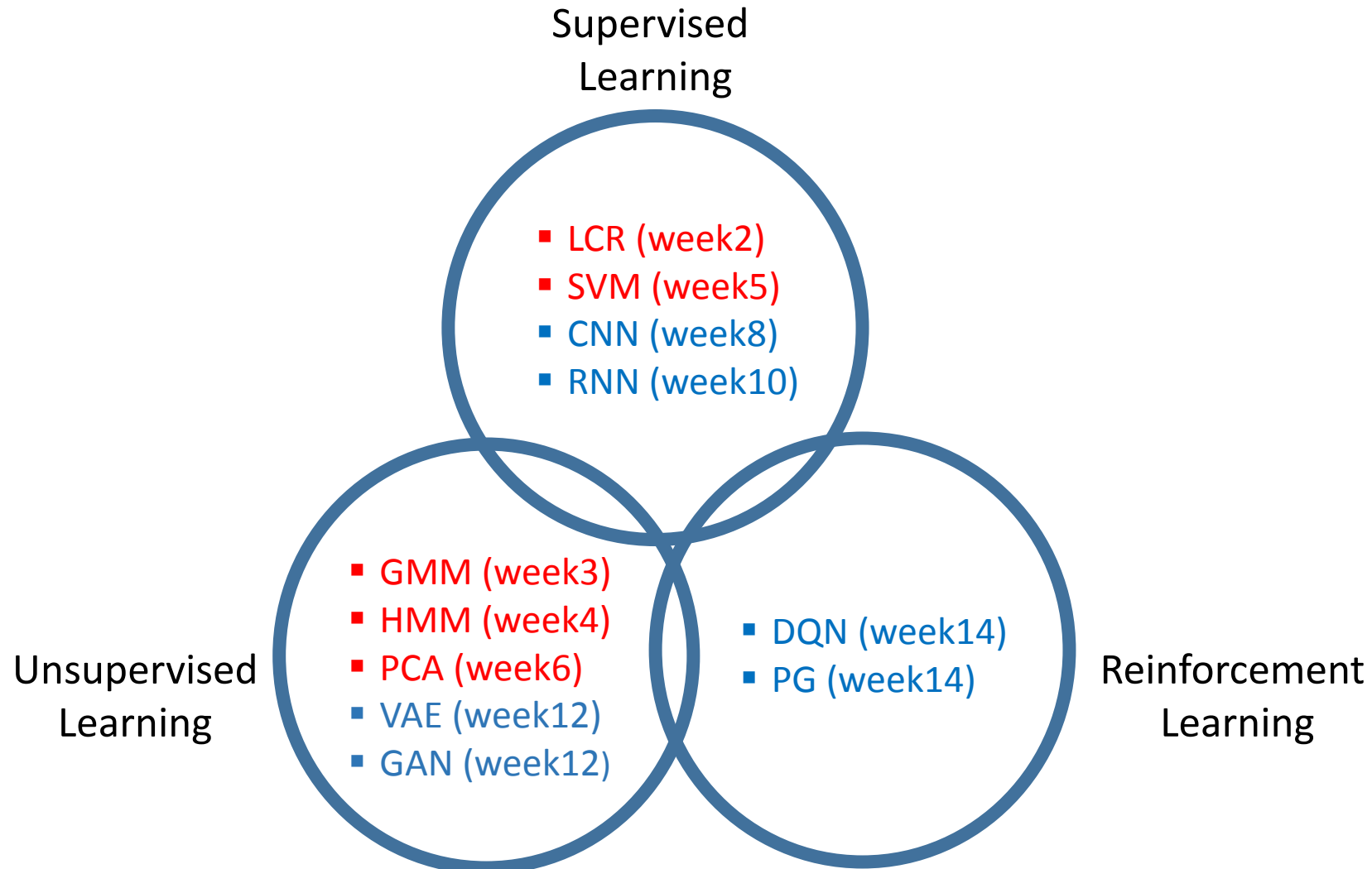
Practical Machine Learning

Lecture 7

Neural Networks

Dr. Suyong Eum

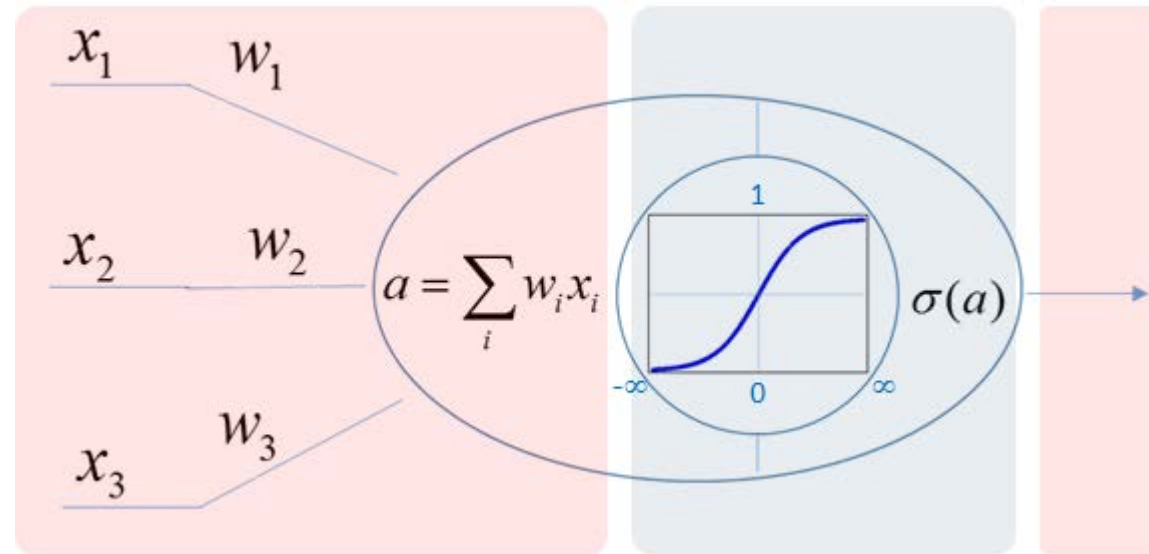
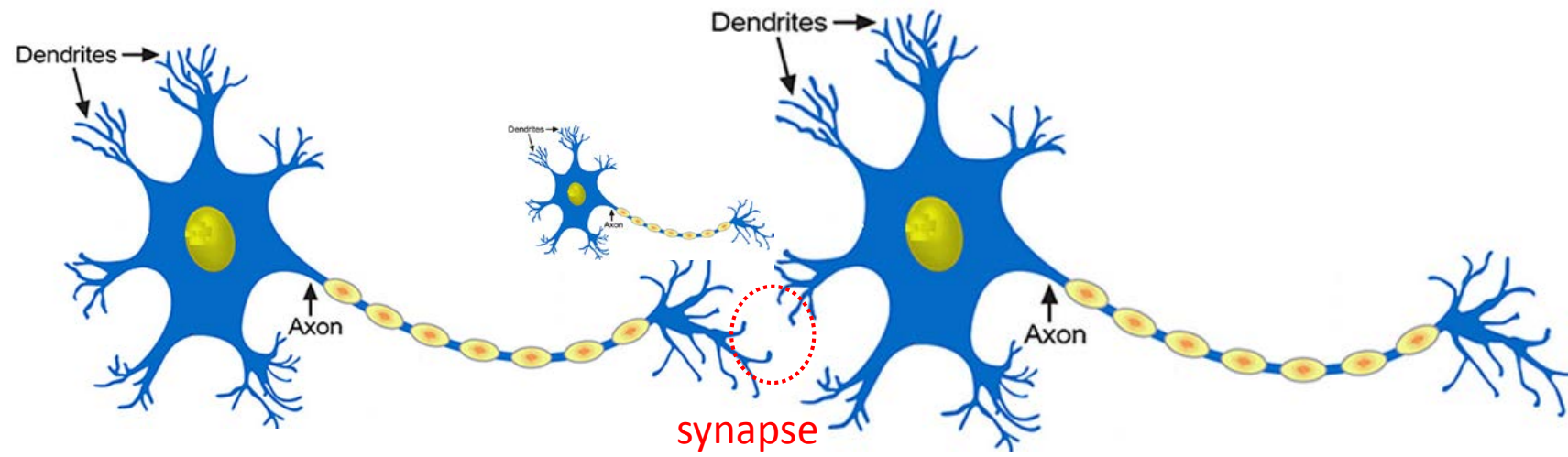




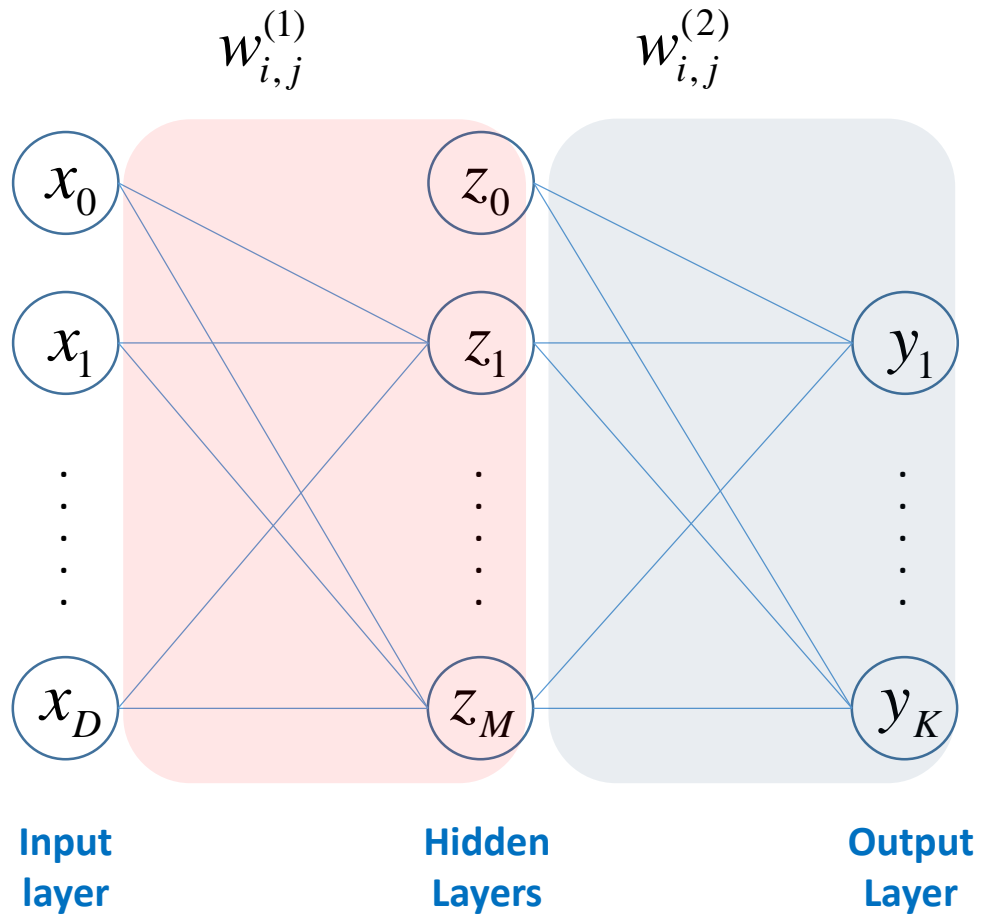
You are going to learn

- ❑ Terminology for neural networks
- ❑ Its basic operation with a multi-layer perceptron model (MLP)
- ❑ How to train a neural network
 - Backpropagation
 - An example of backpropagation

A bio-inspired approach



Terminology in neural networks



$w_{i,j}^{(\ell)}$: weight on a link at layer (ℓ) between node i and j

❑ Two types of layers:

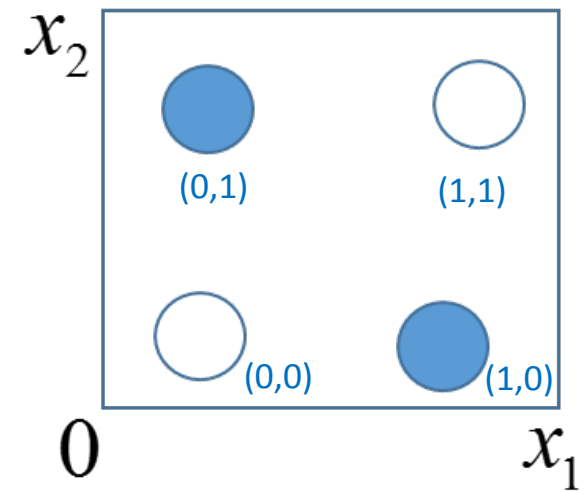
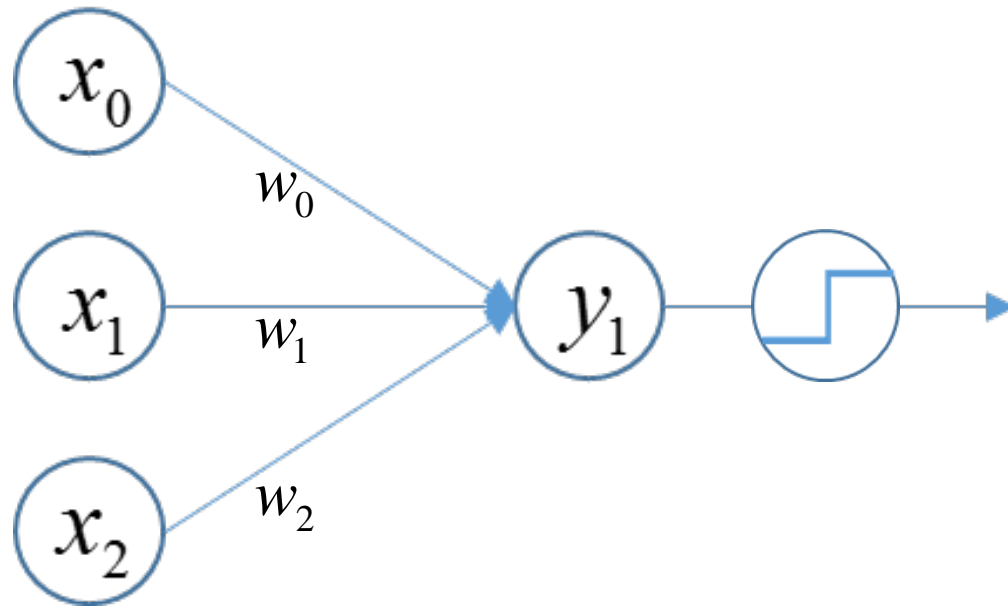
- Nodes: Input/Hidden/Output layer
- Links: Hidden/Output layer

❑ In general, a standard L -layer neural network consists of

- an input layer,
- $(L-1)$ hidden layers,
- an output layer.

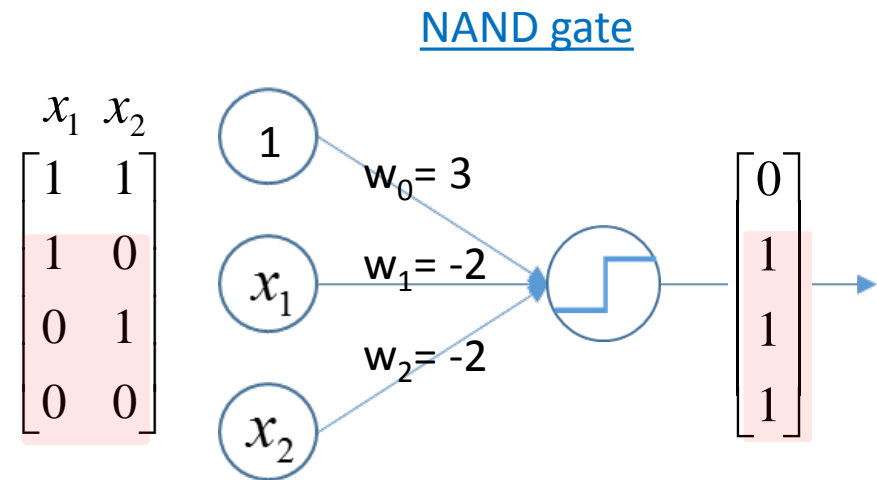
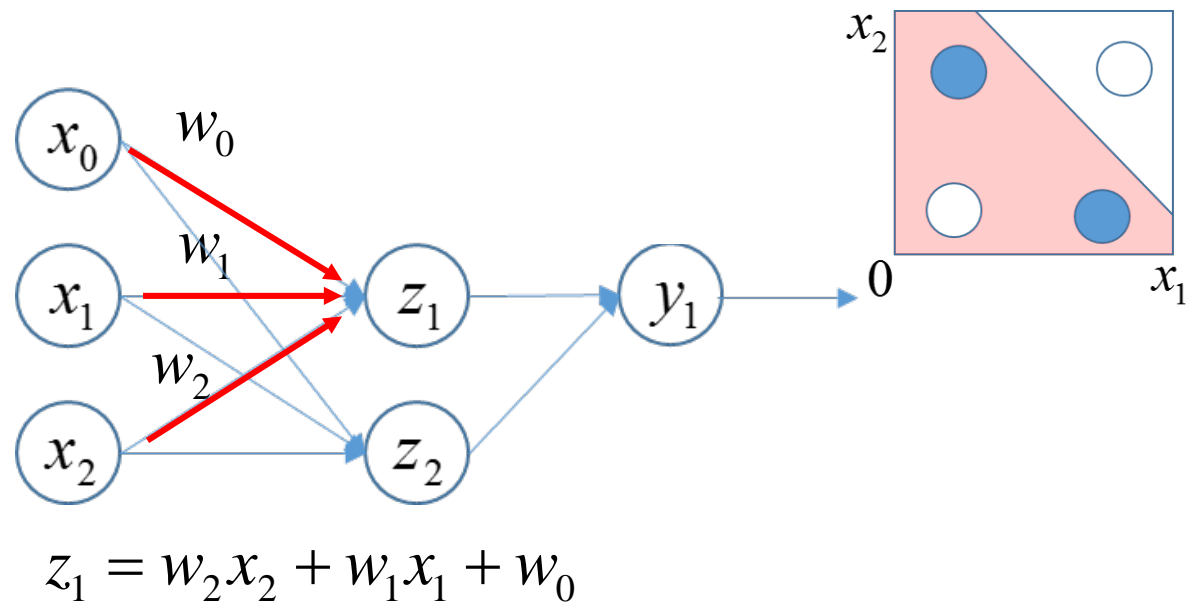
Role of hidden layers

$$y_1 = w_2x_2 + w_1x_1 + w_0$$

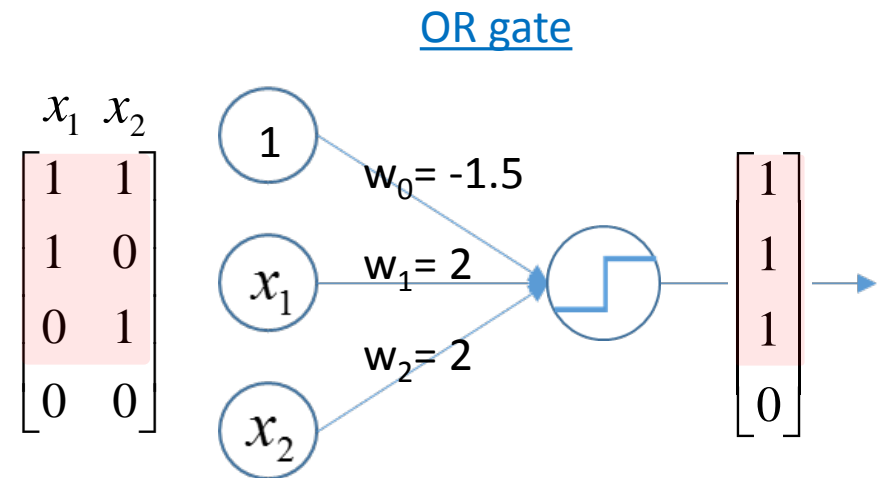
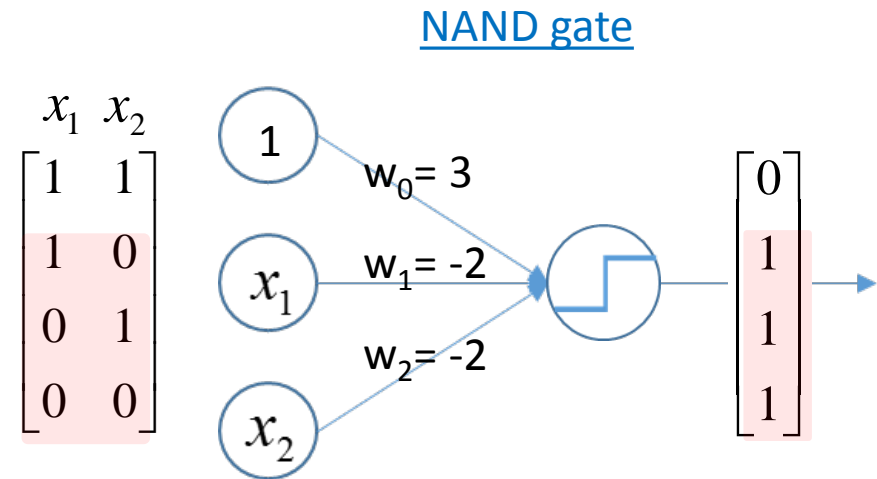
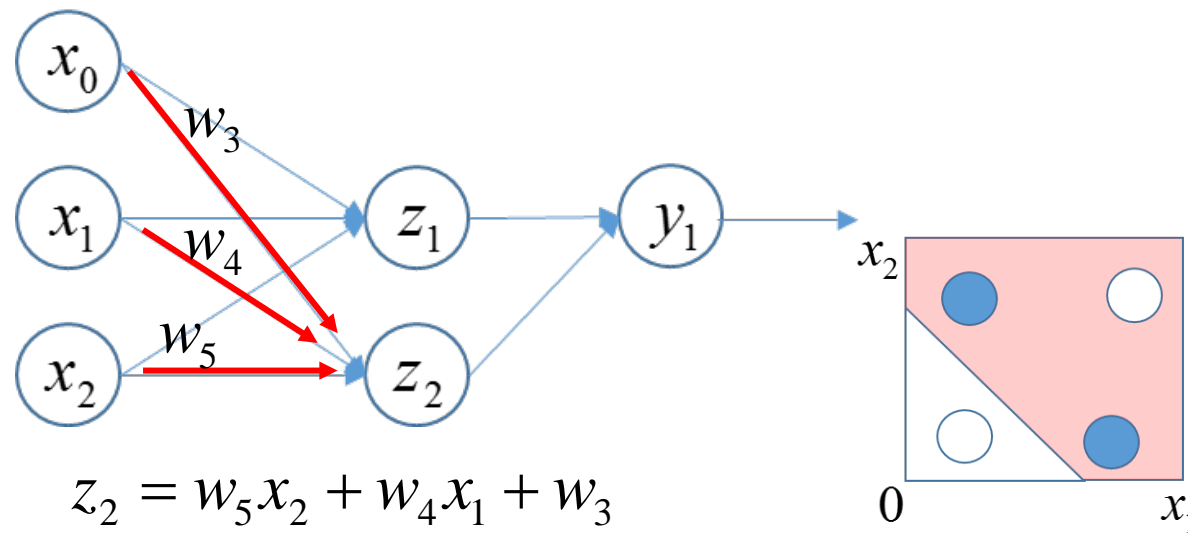
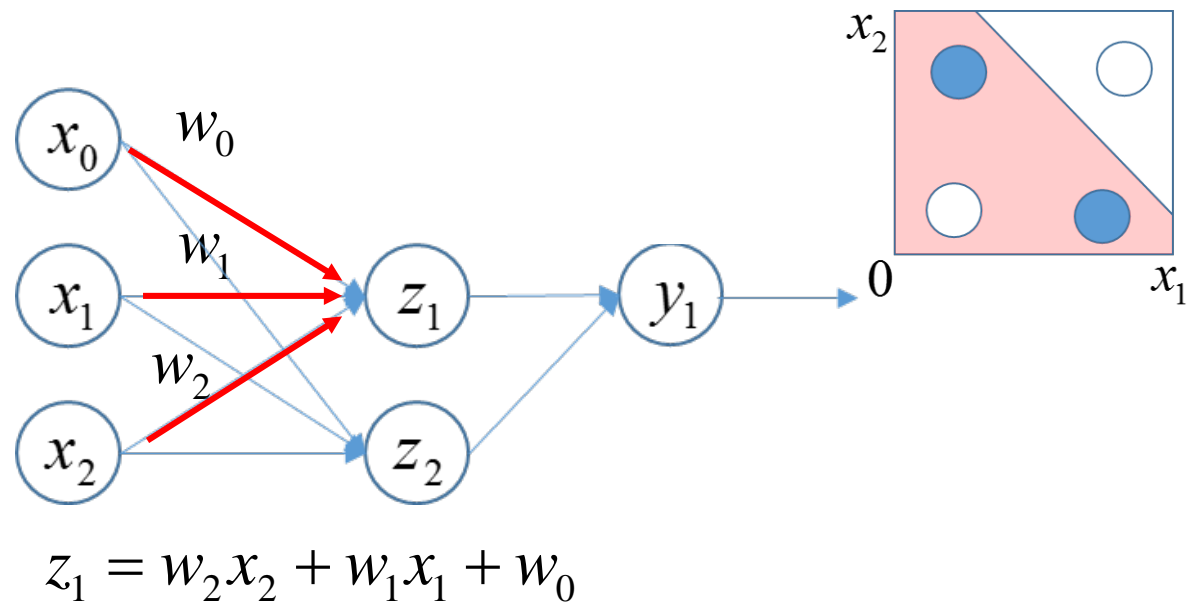


XOR problem

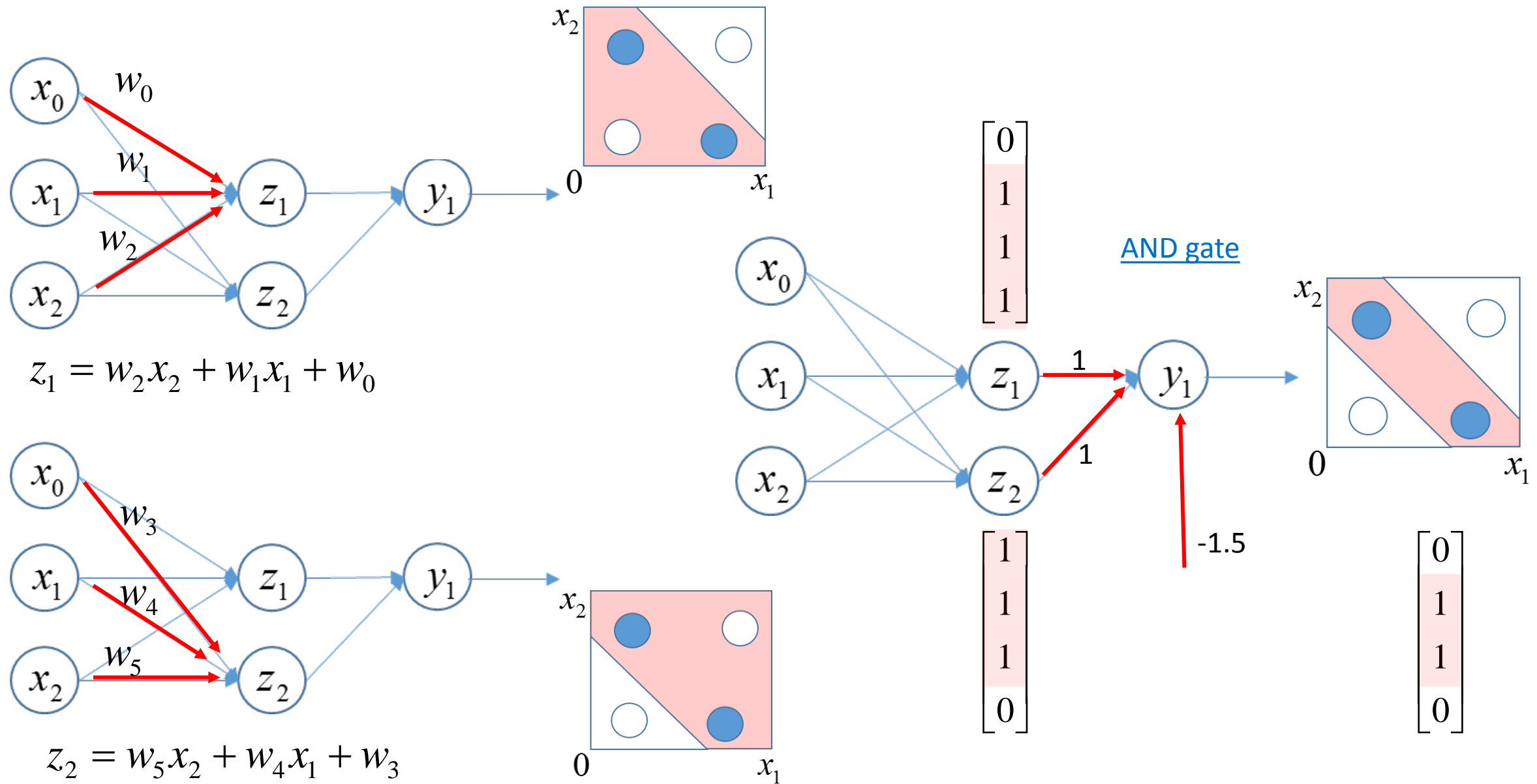
Role of hidden layers: a simple example



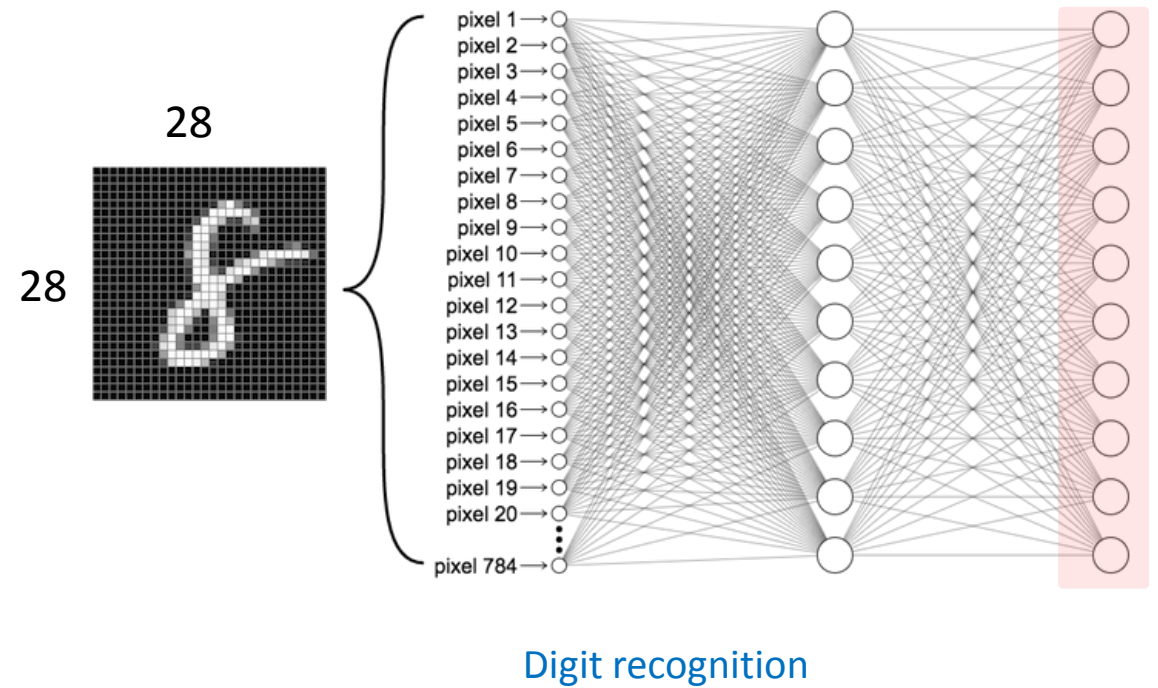
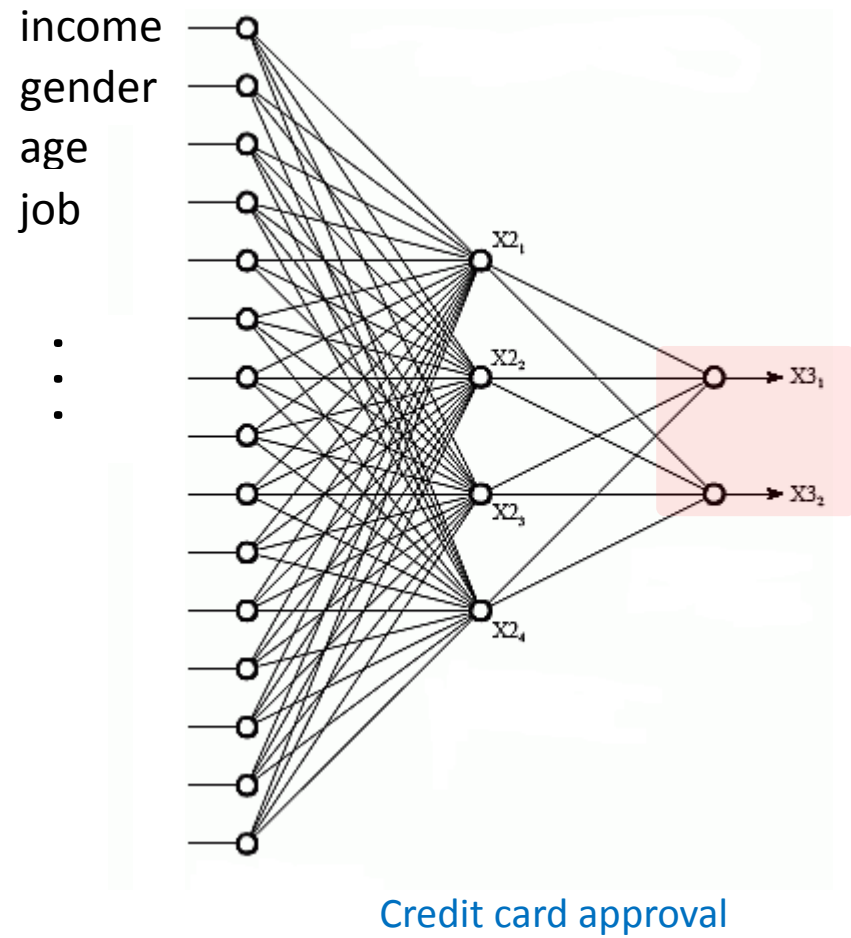
Role of hidden layers: a simple example



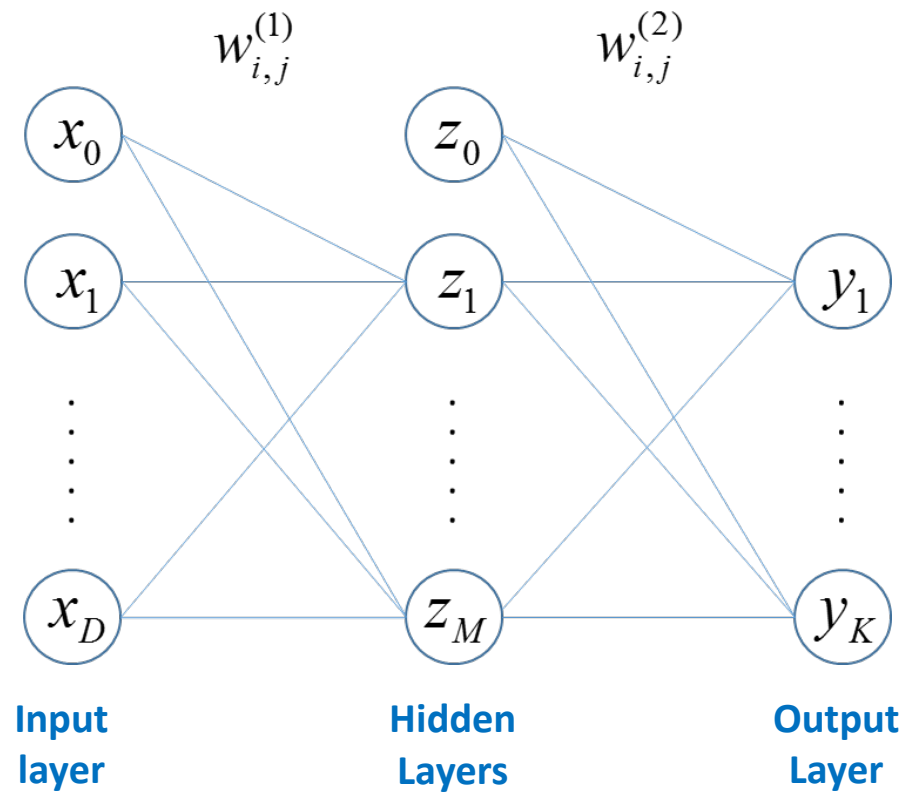
Role of hidden layers: a simple example



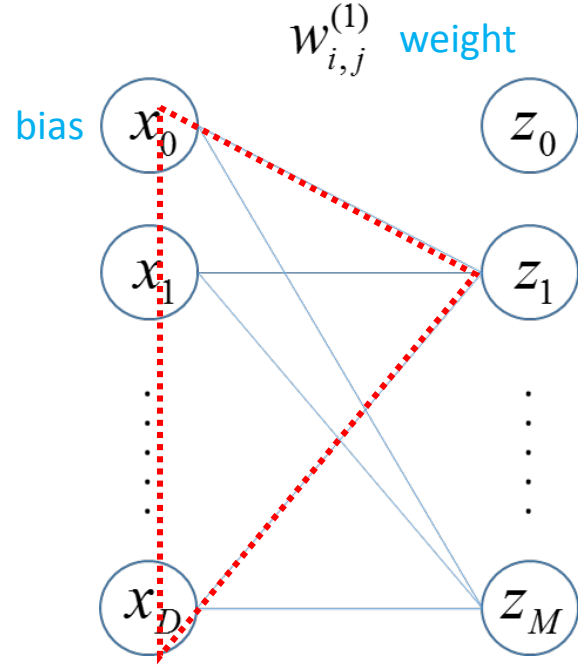
Example structures of neural networks



A neural network model



A neural network model: bias and weight



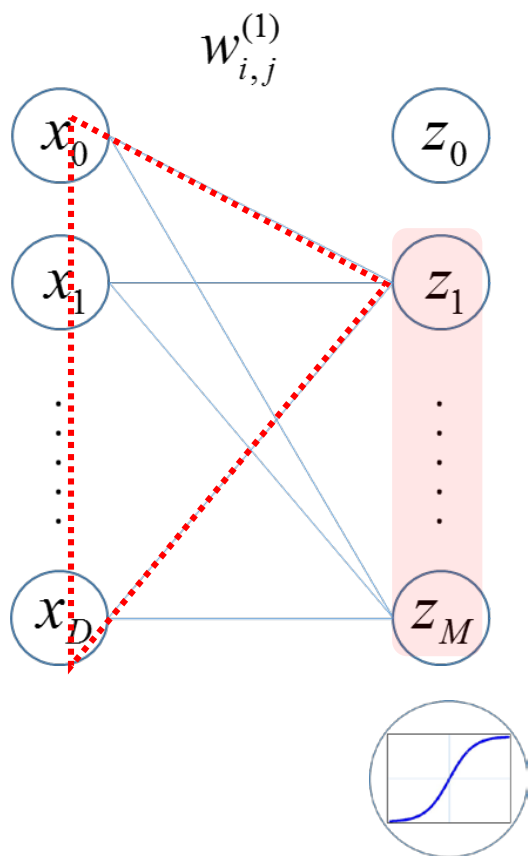
- Fully connected vs Partially connected
- Bias and weight initialization
 - Xavier (2010) / He (2015)

$$w_{1,1}^{(1)}x_1 + w_{2,1}^{(1)}x_2 + w_{0,1}^{(1)}x_0$$

bias

$$x_0 = 1$$

A neural network model: activation function



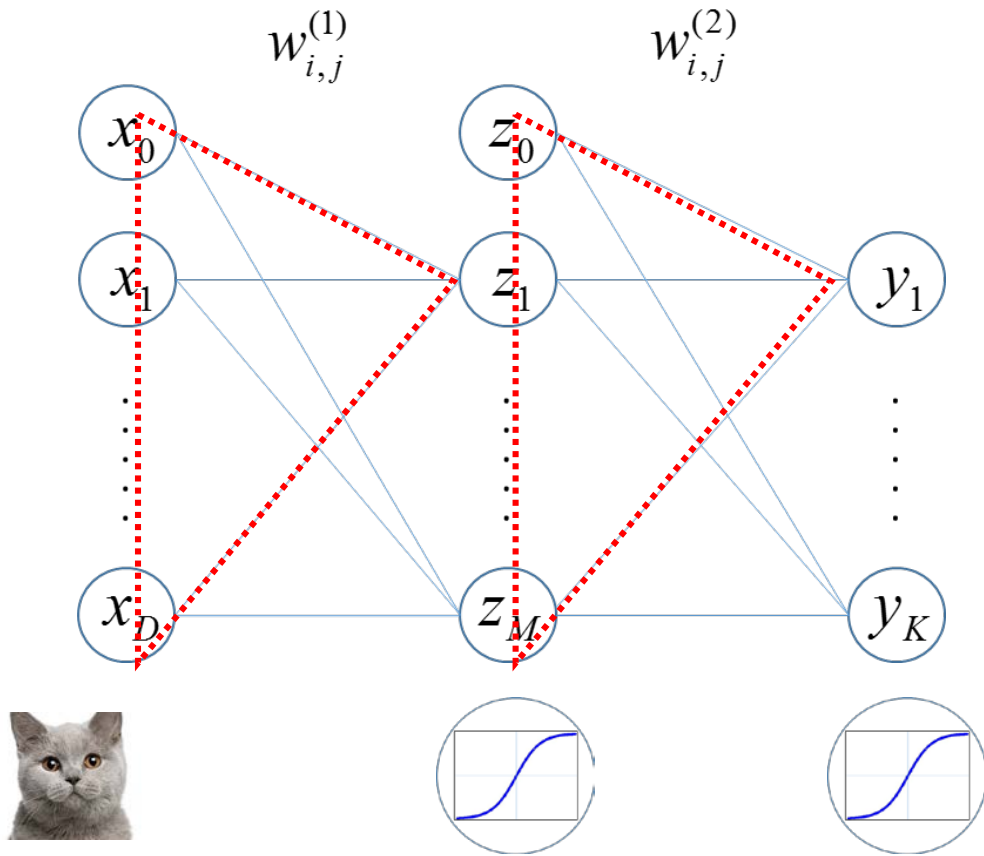
$$w_{1,1}^{(1)} x_1 + w_{2,1}^{(1)} x_2 + w_{0,1}^{(1)} x_0$$

Identity		$f(x) = x$	$f'(x) = 1$
Binary step		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x \neq 0 \\ ? & \text{for } x = 0 \end{cases}$
Logistic (a.k.a. Sigmoid or Soft step)		$f(x) = \sigma(x) = \frac{1}{1 + e^{-x}}$ ^[1]	$f'(x) = f(x)(1 - f(x))$
TanH		$f(x) = \tanh(x) = \frac{(e^x - e^{-x})}{(e^x + e^{-x})}$	$f'(x) = 1 - f(x)^2$
Rectified linear unit (ReLU) ^[11]		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
Leaky rectified linear unit (Leaky ReLU) ^[12]		$f(x) = \begin{cases} 0.01x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0.01 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$

❑ Vanishing gradient problem

❑ Convergence speed (6 times faster) [\[http://www.cs.toronto.edu/~fritz/absps/imagenet.pdf\]](http://www.cs.toronto.edu/~fritz/absps/imagenet.pdf)

A neural network model: cross entropy with softmax



$$w_{1,1}^{(1)} x_1 + w_{2,1}^{(1)} x_2 + w_{0,1}^{(1)} x_0$$

$$w_{1,1}^{(2)} z_1 + w_{2,1}^{(2)} z_2 + w_{0,1}^{(2)} z_0$$

output
-5
-1
1
5

frog
bird
dog
cat

Sigmoid
0.00669
0.26894
0.73106
0.99331

Normalization
0.00334
0.13447
0.36553
0.49666

y Softmax
0.00004
0.00243
0.01794
0.97959

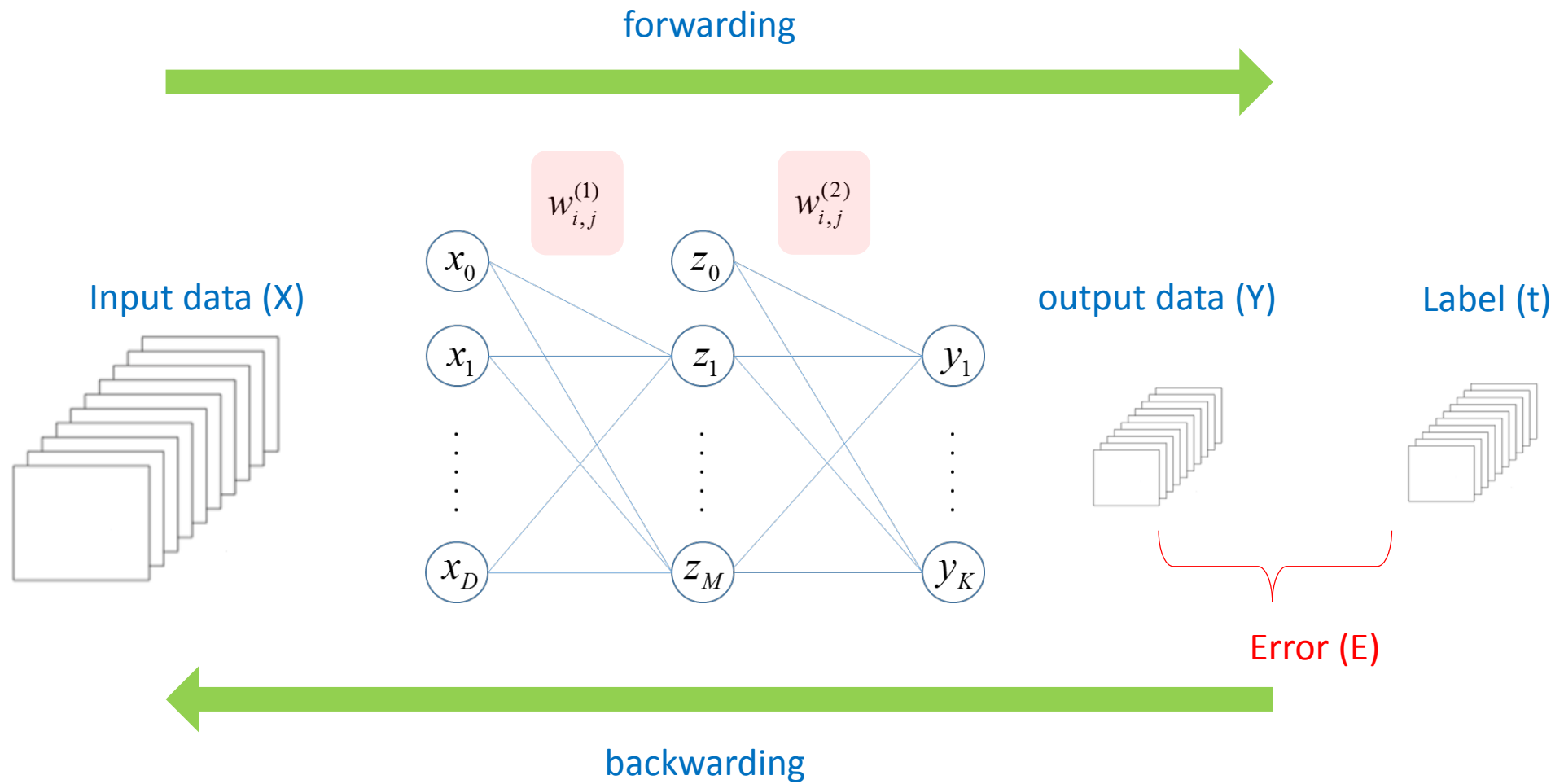
t

Label
0
0
0
1

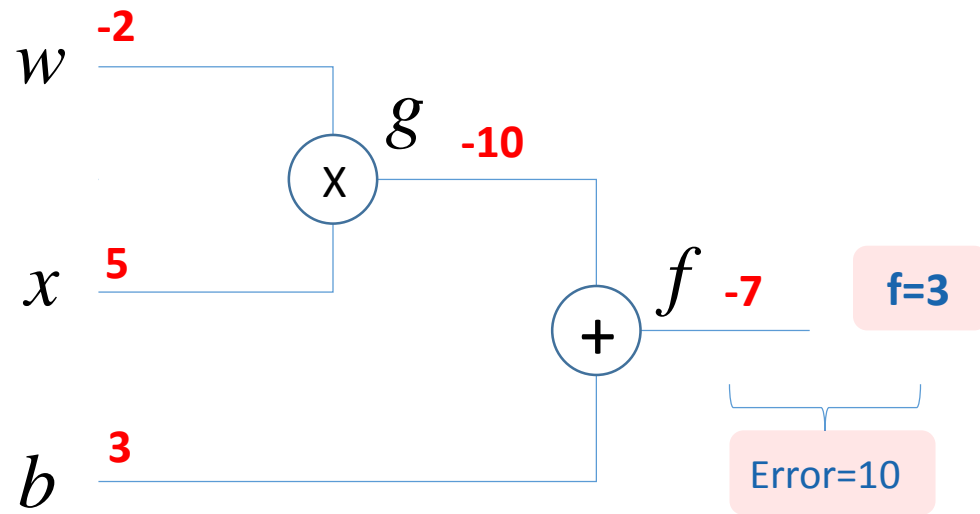
$$H(y) = - \sum_i t_i \log(y_i) = 0.020621$$

Operation

Overview of the operation



Backpropagation: a toy example



$$\frac{\partial f}{\partial w} = \frac{\partial f}{\partial g} \frac{\partial g}{\partial w} = x = 5$$

$$\frac{\partial f}{\partial b} = 1$$

$$f = g + b$$

$$g = wx$$

- ❑ Assuming that the value of f should be “3”.
- ❑ How to update variables which you are interested?

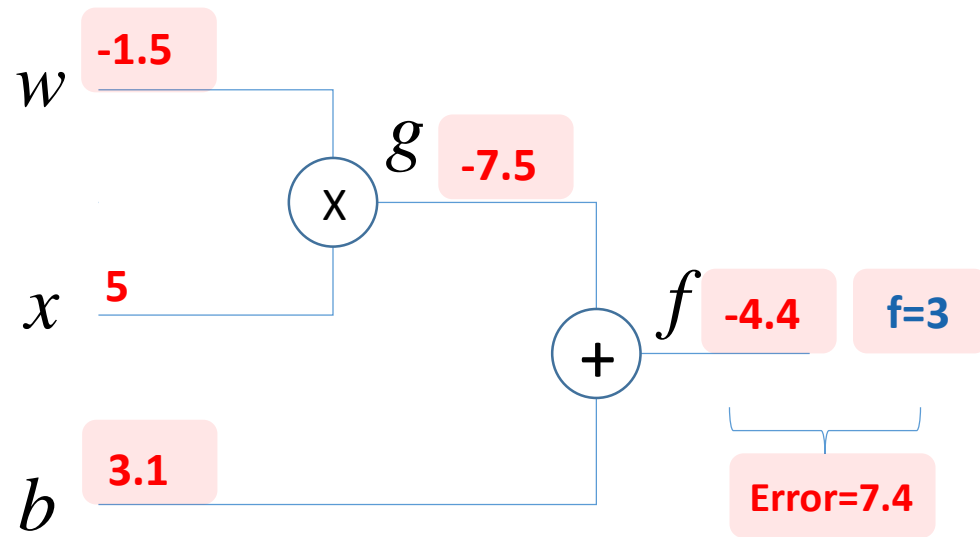
$$W_{new} = W_{old} + \eta \frac{\partial f}{\partial w_i}$$

$$b_{new} = b_{old} + \eta \frac{\partial f}{\partial b}$$

$$W_{new} = -2 + 0.1 \times 5 = -1.5$$

$$b_{new} = 3 + 0.1 \times 1 = 3.1$$

Backpropagation: a toy example



$$\frac{\partial f}{\partial w} = \frac{\partial f}{\partial g} \frac{\partial g}{\partial w} = x = 5$$

$$\frac{\partial f}{\partial b} = 1$$

$$f = g + b$$

$$g = wx$$

- ❑ Assuming that the value of f should be “3”.
- ❑ How to update variables which you are interested?

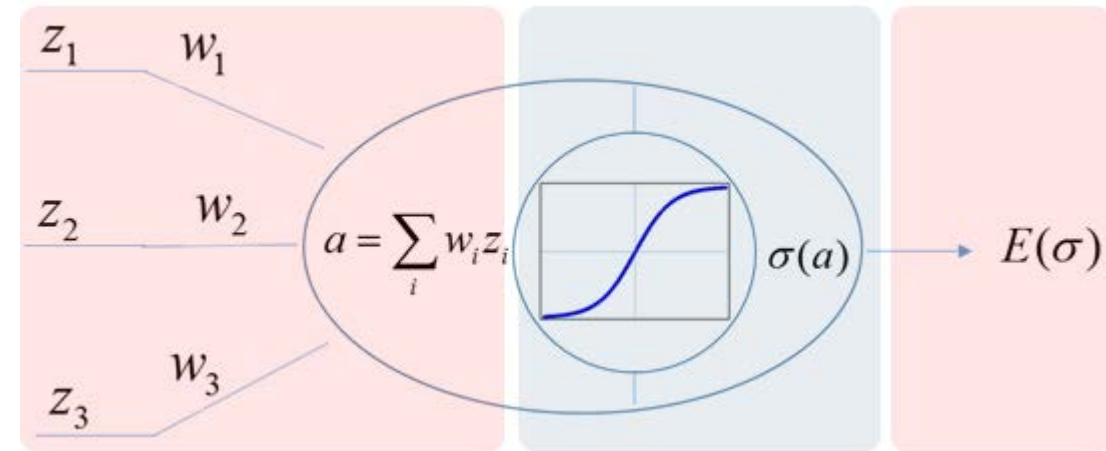
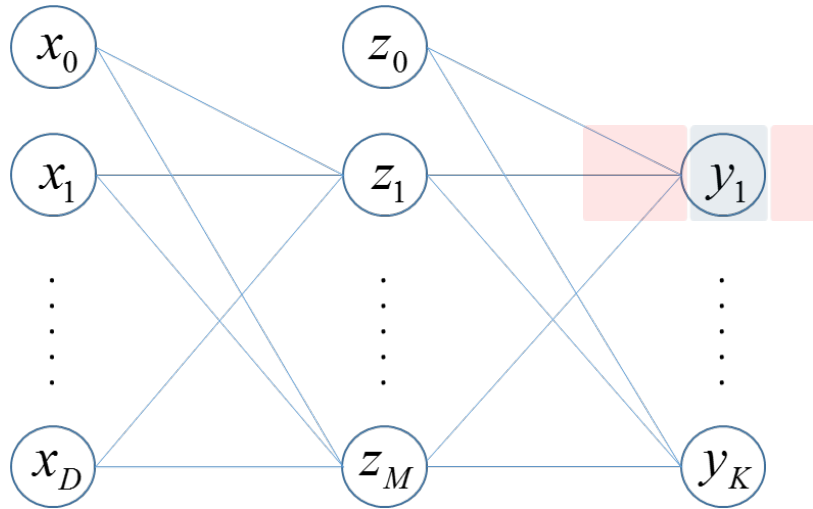
$$W_{new} = W_{old} + \eta \frac{\partial f}{\partial w_i}$$

$$b_{new} = b_{old} + \eta \frac{\partial f}{\partial b}$$

$$W_{new} = -2 + 0.1 \times 5 = -1.5$$

$$b_{new} = 3 + 0.1 \times 1 = 3.1$$

Backpropagation in a neuron



Linear summation
function

Activation
function

Error
function

$$\frac{\partial E(\sigma(a))}{\partial w} = \frac{\partial E(\sigma(a))}{\partial \sigma} \times \frac{\partial \sigma}{\partial a} \times \frac{\partial a}{\partial w}$$

$$a = w_1 z_1 + w_2 z_2 + w_3 z_3$$

$$\sigma(a) = \frac{1}{1 + e^{-a}}$$

$$E(\sigma(a)) = \frac{1}{2} \sum_{i=1} (t_i - \sigma(a))^2$$

Backpropagation in a neuron

Linear summation function

$$a = w_1 z_1 + w_2 z_2 + w_3 z_3$$

$$\partial a / \partial w_1 = z_1$$

$$\partial a / \partial w_2 = z_2$$

$$\partial a / \partial w_3 = z_3$$

Activation function

$$\sigma(a) = \frac{1}{1 + e^{-a}}$$

$$\frac{\partial \sigma}{\partial a} = \sigma(a)(1 - \sigma(a))$$

Error function

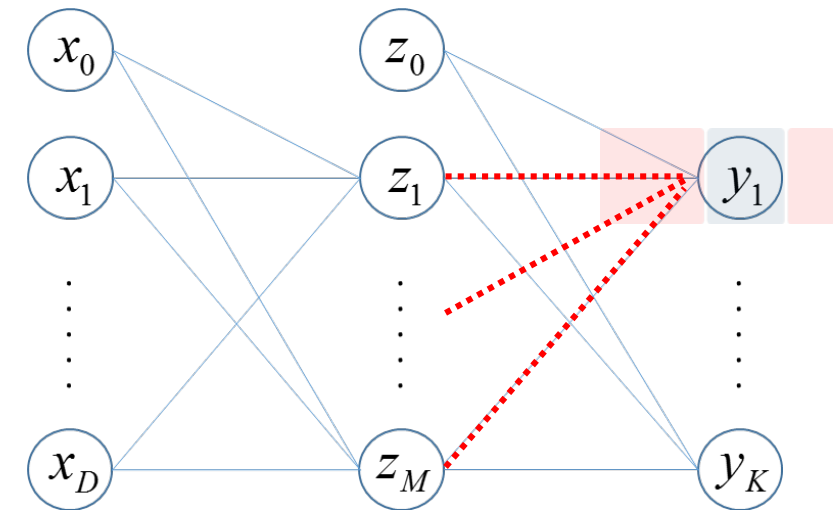
$$E(\sigma(a)) = \frac{1}{2} \sum_{i=1} (t_i - \sigma(a))^2$$

$$\frac{\partial E(\sigma(a))}{\partial \sigma} = -(t_i - \sigma(a))$$

$$\frac{\partial E(w)}{\partial w_1} = \frac{\partial E(\sigma(a))}{\partial \sigma} \times \frac{\partial \sigma}{\partial a} \times \frac{\partial a}{\partial w_1}$$

$$\frac{\partial E(w)}{\partial w_2} = \frac{\partial E(\sigma(a))}{\partial \sigma} \times \frac{\partial \sigma}{\partial a} \times \frac{\partial a}{\partial w_2}$$

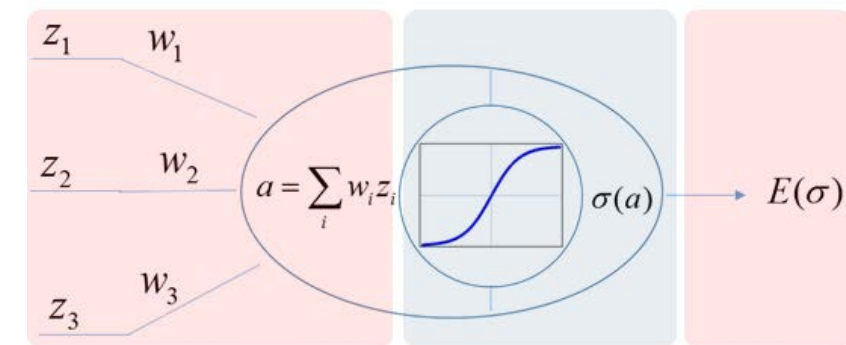
$$\frac{\partial E(w)}{\partial w_3} = \frac{\partial E(\sigma(a))}{\partial \sigma} \times \frac{\partial \sigma}{\partial a} \times \frac{\partial a}{\partial w_3}$$



Linear summation function

Activation function

Error function



Backpropagation in the next layer

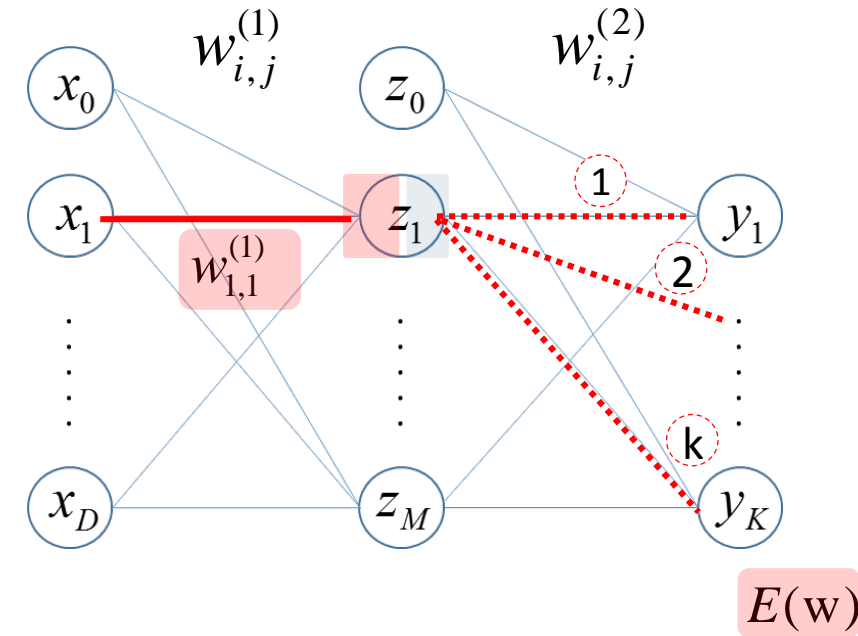
$$\frac{\partial E(\mathbf{w})}{\partial w_{1,1}^{(1)}} = \frac{\partial E(\mathbf{w})}{\partial z_{1(\text{output})}} \times \frac{\partial z_{1(\text{output})}}{\partial z_{1(\text{input})}} \times \frac{\partial z_{1(\text{input})}}{\partial w_{1,1}^{(1)}}$$

$$\frac{\partial E(\mathbf{w})}{\partial z_{1(\text{output})}} = \overset{1}{\frac{\partial E(\mathbf{w})_{y_1}}{\partial z_{1(\text{output})}}} + \overset{2}{\frac{\partial E(\mathbf{w})_{y_2}}{\partial z_{1(\text{output})}}} + \dots + \overset{k}{\frac{\partial E(\mathbf{w})_{y_K}}{\partial z_{1(\text{output})}}}$$

$$\frac{\partial E(\mathbf{w})_{y_1}}{\partial z_{1(\text{output})}} = \frac{\partial E(\mathbf{w})_{y_1}}{\partial y_{1(\text{output})}} \times \frac{\partial y_{1(\text{output})}}{\partial y_{1(\text{input})}} \times \frac{\partial y_{1(\text{input})}}{\partial z_{1(\text{output})}}$$

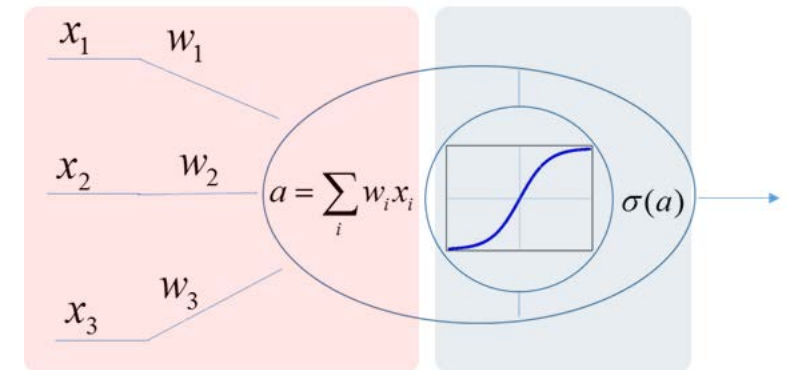
Previously obtained $w_{1,1}^{(2)}$

$$y_1 = w_{1,1}^{(2)} z_1 + w_{2,1}^{(2)} z_2 + w_{0,1}^{(2)} z_0$$



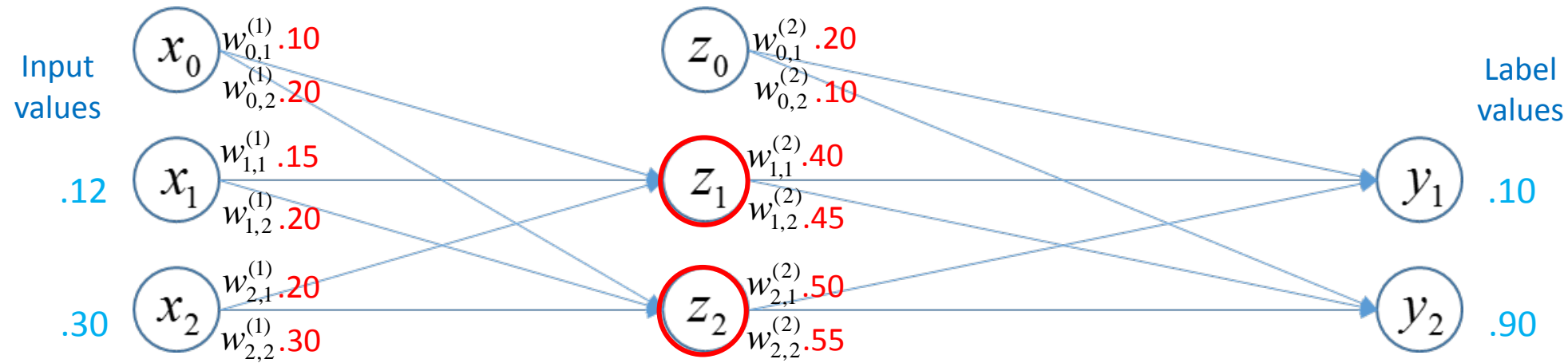
Linear summation
function

Activation
function

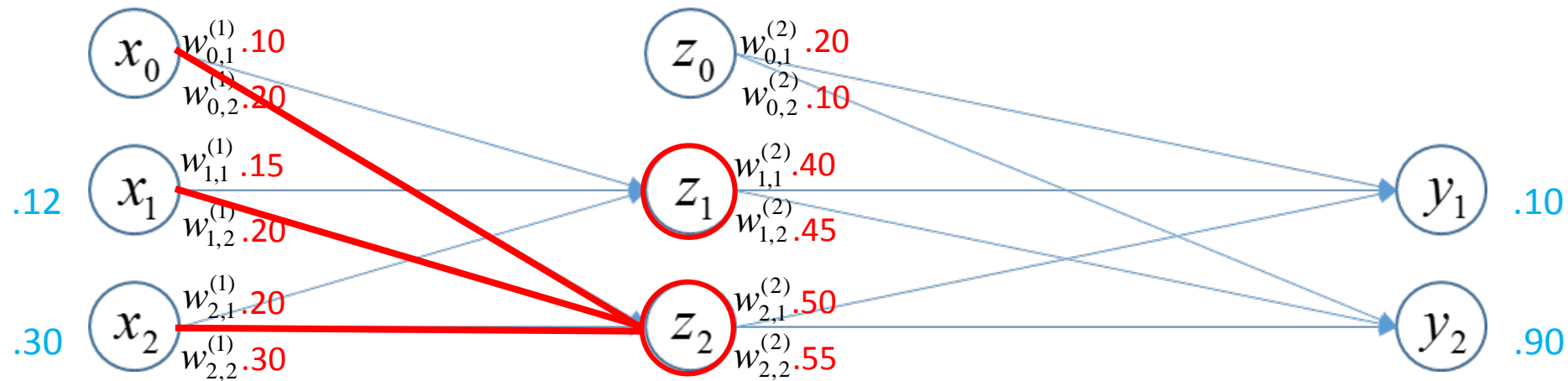


Backpropagation example

Backpropagation algorithm - Forwarding



Backpropagation algorithm - Forwarding



$$a_1 = w_{1,1}^{(1)} x_1 + w_{2,1}^{(1)} x_2 + w_{0,1}^{(1)}$$

$$a_1 = 0.15 \times 0.12 + 0.2 \times 0.3 + 0.1$$

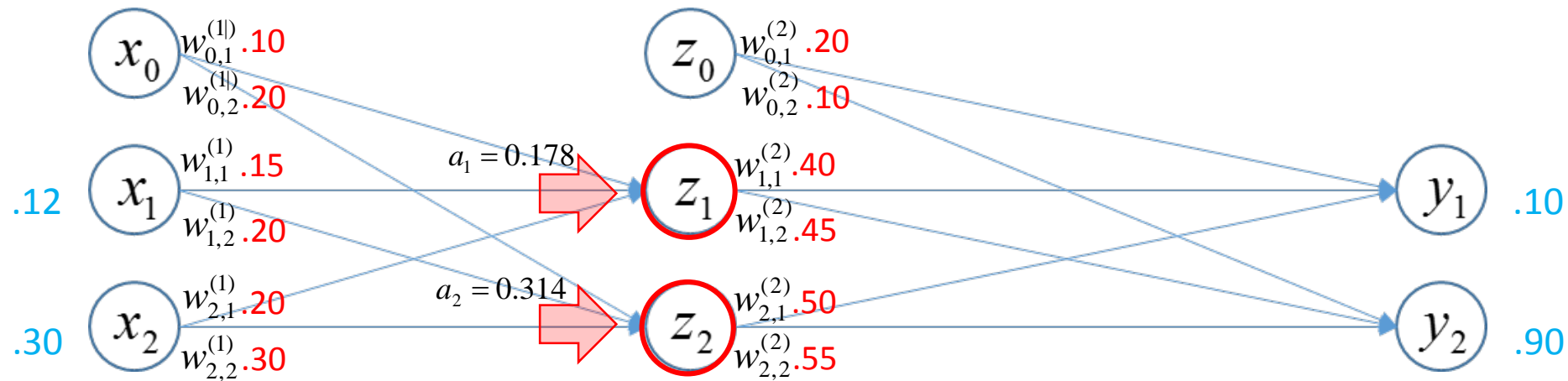
$$a_1 = 0.178$$

$$a_2 = w_{1,2}^{(1)} x_1 + w_{2,2}^{(1)} x_2 + w_{0,2}^{(1)}$$

$$a_2 = 0.2 \times 0.12 + 0.3 \times 0.3 + 0.2$$

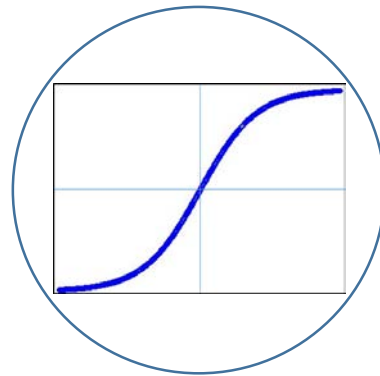
$$a_2 = 0.314$$

Backpropagation algorithm - Forwarding



$$a_1 = 0.178$$

$$a_2 = 0.314$$

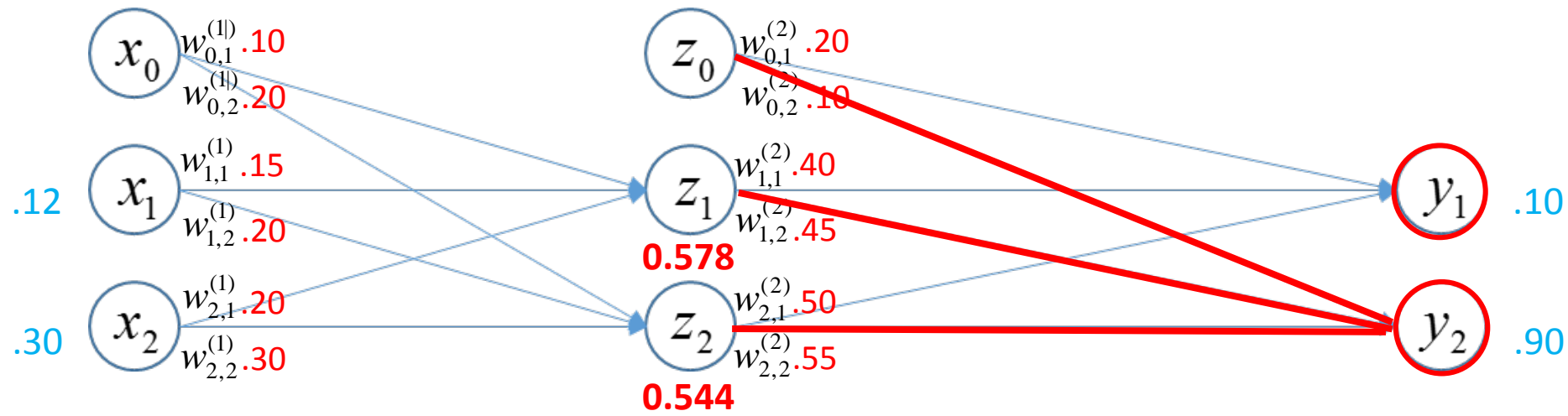


$$\sigma(a_1) = 0.544$$

$$\sigma(a_2) = 0.578$$

$$\sigma(a) = \frac{1}{1 + e^{-a}}$$

Backpropagation algorithm - Forwarding



$$a_1 = w_{1,1}^{(2)} z_1 + w_{2,1}^{(2)} z_2 + w_{0,1}^{(2)}$$

$$a_1 = 0.40 \times 0.578 + 0.5 \times 0.544 + 0.2$$

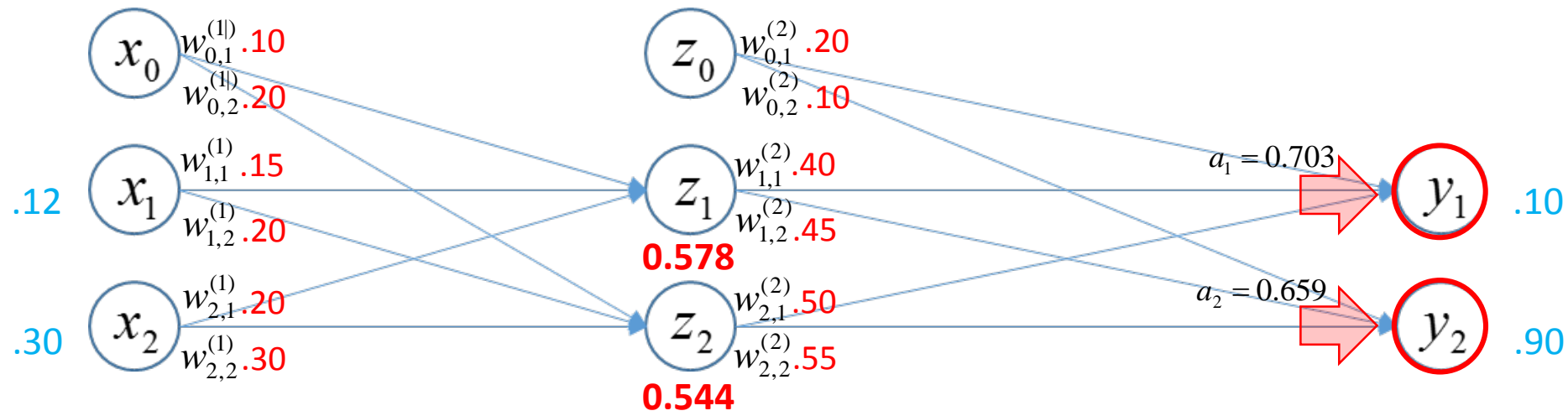
$$a_1 = 0.703$$

$$a_2 = w_{1,2}^{(2)} z_1 + w_{2,2}^{(2)} z_2 + w_{0,2}^{(2)}$$

$$a_2 = 0.45 \times 0.578 + 0.55 \times 0.544 + 0.1$$

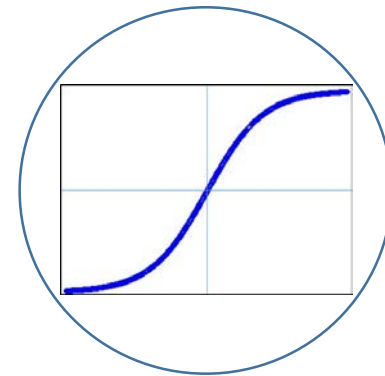
$$a_2 = 0.659$$

Backpropagation algorithm - Forwarding



$$a_1 = 0.703$$

$$a_2 = 0.659$$

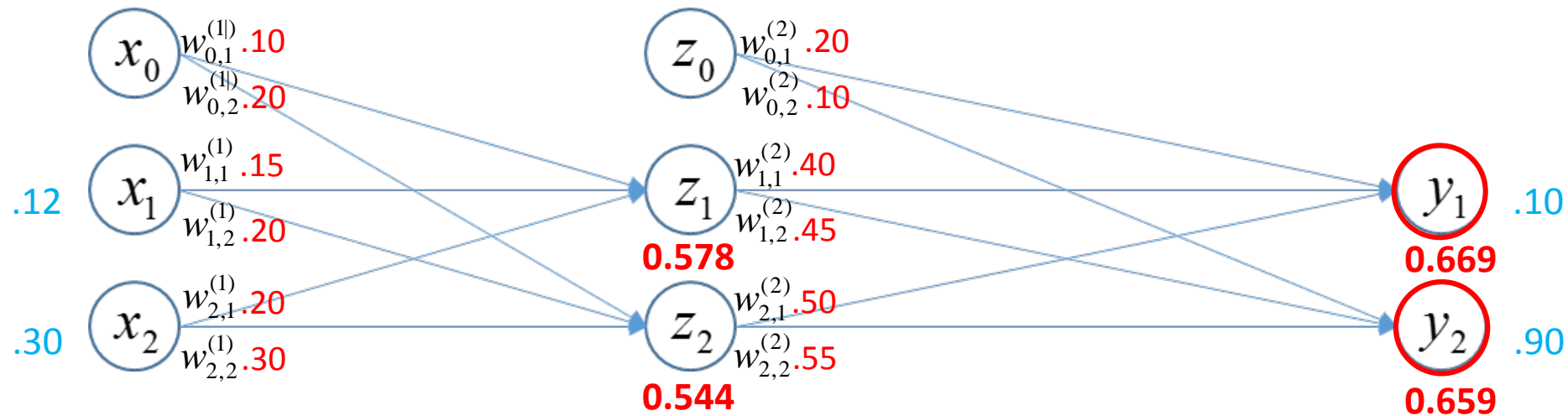


$$\sigma(a_1) = 0.669$$

$$\sigma(a_2) = 0.659$$

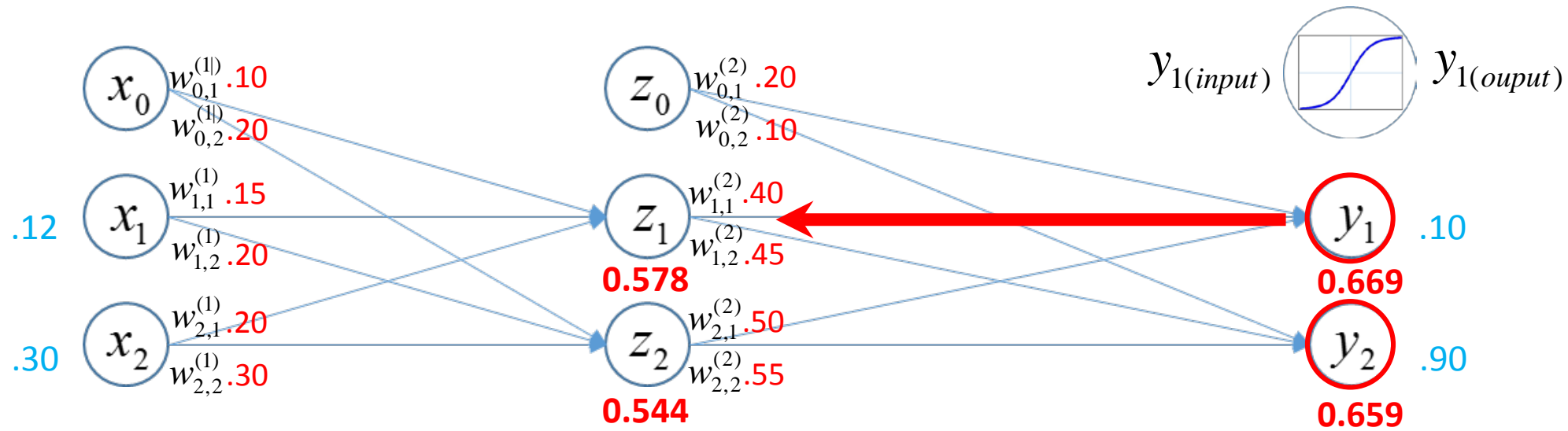
$$\sigma(a) = \frac{1}{1 + e^{-a}}$$

Backpropagation algorithm - Forwarding



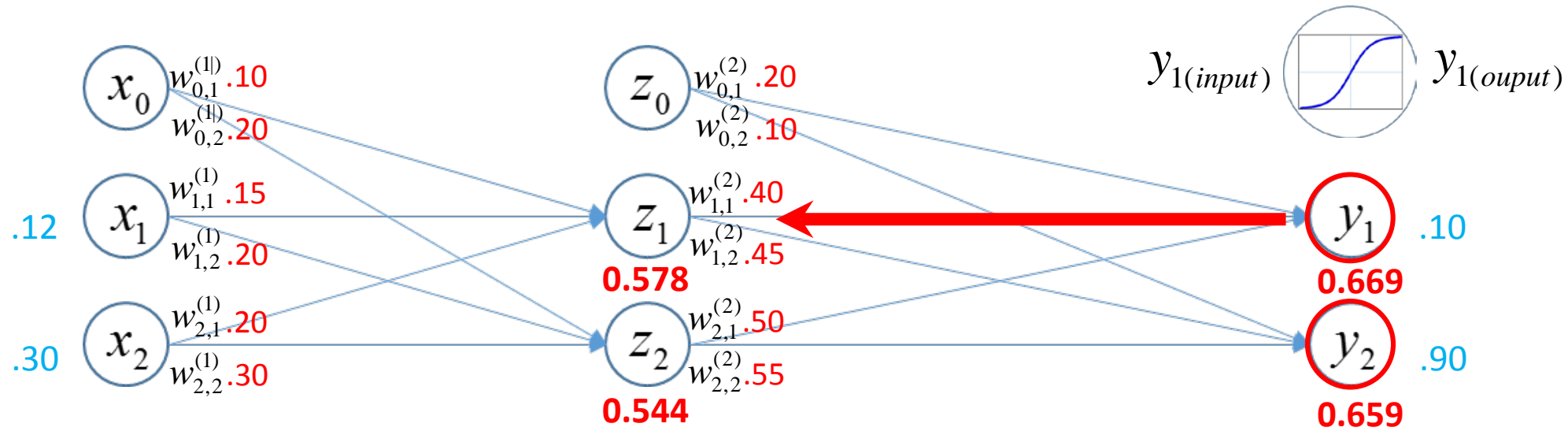
$$E(w) = \frac{1}{2} \left((0.1 - 0.669)^2 + (0.9 - 0.659)^2 \right) = 0.191$$

Backpropagation algorithm - Backwarding



$$\frac{\partial y_{1(input)}}{\partial w_{1,1}^{(2)}} \times \frac{\partial y_{1(output)}}{\partial y_{1(input)}} \times \frac{\partial E(w)}{\partial y_{1(output)}} = \frac{\partial E(w)}{\partial w_{1,1}^{(2)}}$$

Backpropagation algorithm - Backwarding



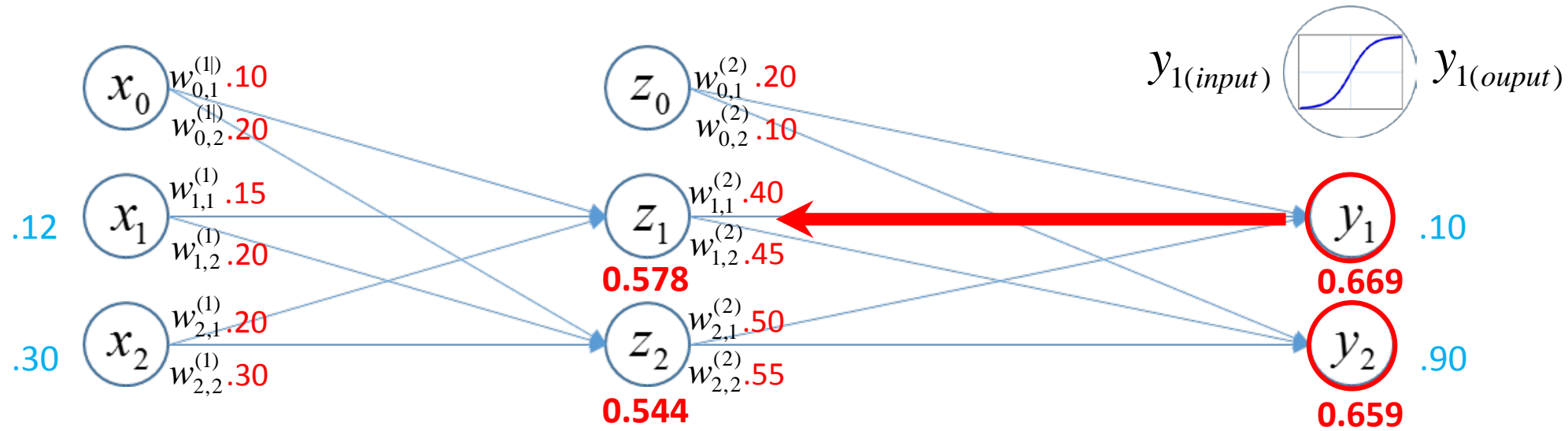
$$\frac{\partial y_{1(input)}}{\partial w_{1,1}^{(2)}} \times \frac{\partial y_{1(output)}}{\partial y_{1(input)}} \times \frac{\partial E(w)}{\partial y_{1(output)}} = \frac{\partial E(w)}{\partial w_{1,1}^{(2)}}$$

0.569

$$E(w) = \frac{1}{2} \left((0.1 - y_{1(output)})^2 + (0.9 - y_{2(output)})^2 \right)$$

$$\frac{\partial E(w)}{\partial y_{1(output)}} = -(0.1 - y_{1(output)}) = -(0.1 - 0.669) = 0.569$$

Backpropagation algorithm - Backwarding



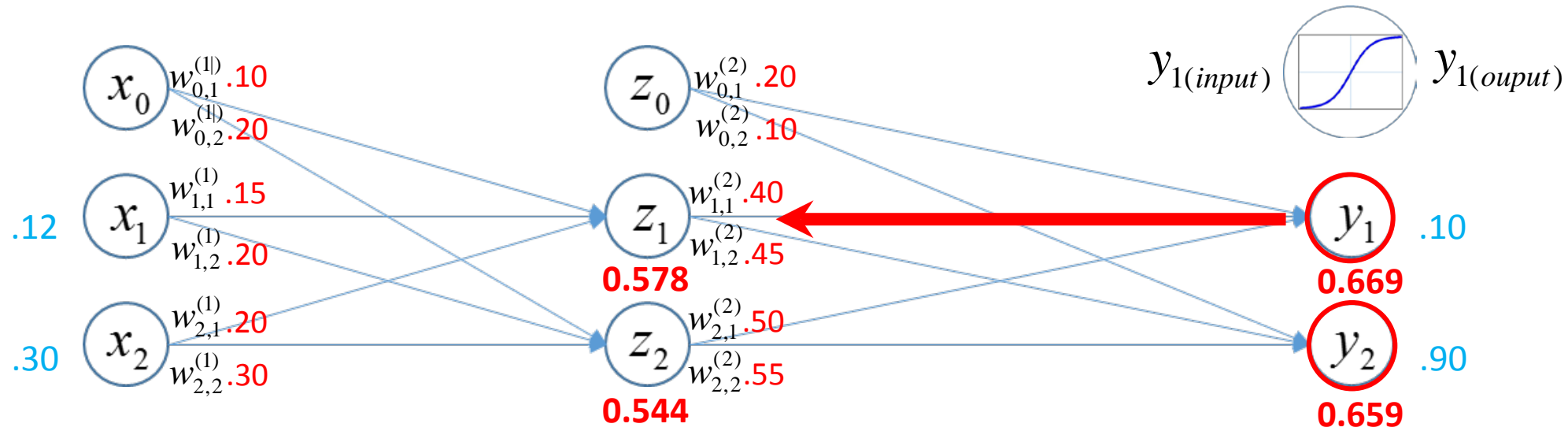
$$y_{1(output)} = \frac{1}{1 + e^{-y_{1(input)}}}$$

$$\begin{aligned} \frac{\partial y_{1(output)}}{\partial y_{1(input)}} &= \sigma(y_{1(input)}) (1 - \sigma(y_{1(input)})) \\ &= y_{1(output)} (1 - y_{1(output)}) \\ &= 0.669 \times (1 - 0.669) = 0.221 \end{aligned}$$

$$\frac{\partial y_{1(input)}}{\partial w_{1,1}^{(2)}} \times \frac{\partial y_{1(output)}}{\partial y_{1(input)}} \times \frac{\partial E(w)}{\partial y_{1(output)}} = \frac{\partial E(w)}{\partial w_{1,1}^{(2)}}$$

0.221 **0.569**

Backpropagation algorithm - Backwarding



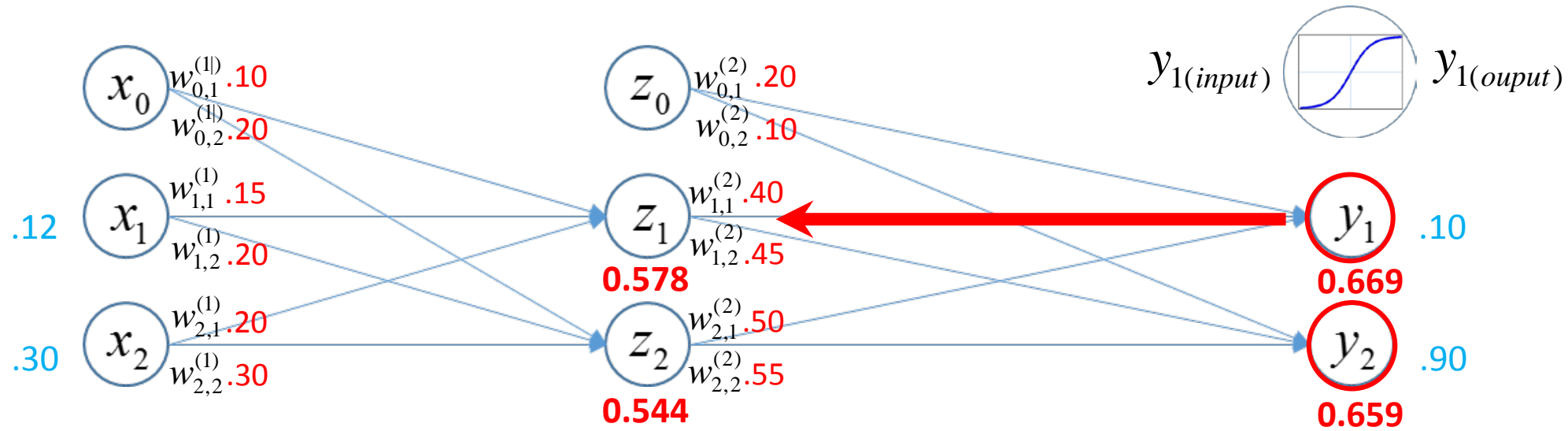
$$y_{1(input)} = w_{1,1}^{(2)} z_1 + w_{2,1}^{(2)} z_2 + w_{0,1}^{(2)}$$

$$\frac{\partial y_{1(input)}}{\partial w_{1,1}^{(2)}} = z_1 = 0.578$$

$$\frac{\partial y_{1(input)}}{\partial w_{1,1}^{(2)}} \times \frac{\partial y_{1(output)}}{\partial y_{1(input)}} \times \frac{\partial E(w)}{\partial y_{1(output)}} = \frac{\partial E(w)}{\partial w_{1,1}^{(2)}}$$

0.578
0.221
0.569

Backpropagation algorithm - Backwarding



$$w_{1,1}^{(2)*} = w_{1,1}^{(2)} + \eta \frac{\partial E(w)}{\partial w_{1,1}^{(2)}}$$

$$= 0.4 + 0.5 \times 0.0727$$

$$= \mathbf{0.436}$$

Same procedures are applied for

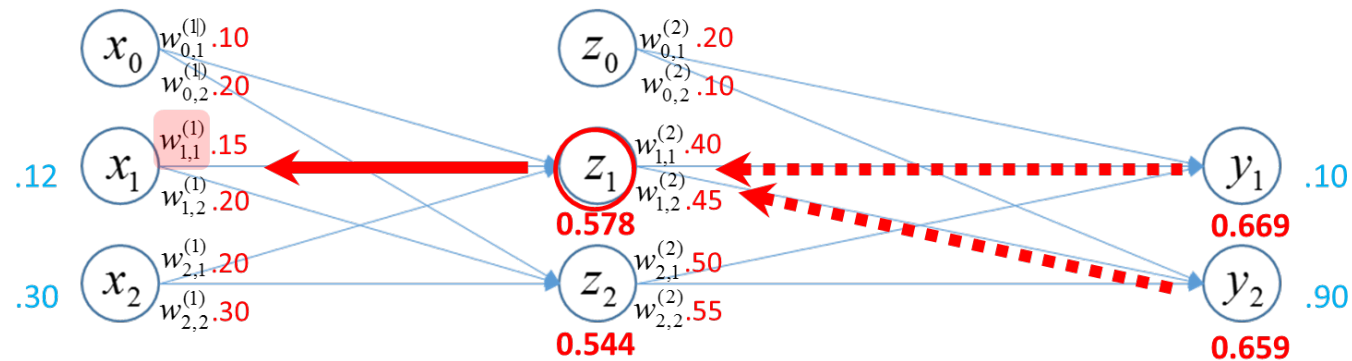
$$w_{1,2}^{(2)}, w_{2,1}^{(2)}, w_{2,2}^{(2)}$$

$$\frac{\partial y_{1(input)}}{\partial w_{1,1}^{(2)}} \times \frac{\partial y_{1(output)}}{\partial y_{1(input)}} \times \frac{\partial E(w)}{\partial y_{1(output)}} = \frac{\partial E(w)}{\partial w_{1,1}^{(2)}}$$

$$\mathbf{0.578} \quad \mathbf{0.221} \quad \mathbf{0.569}$$

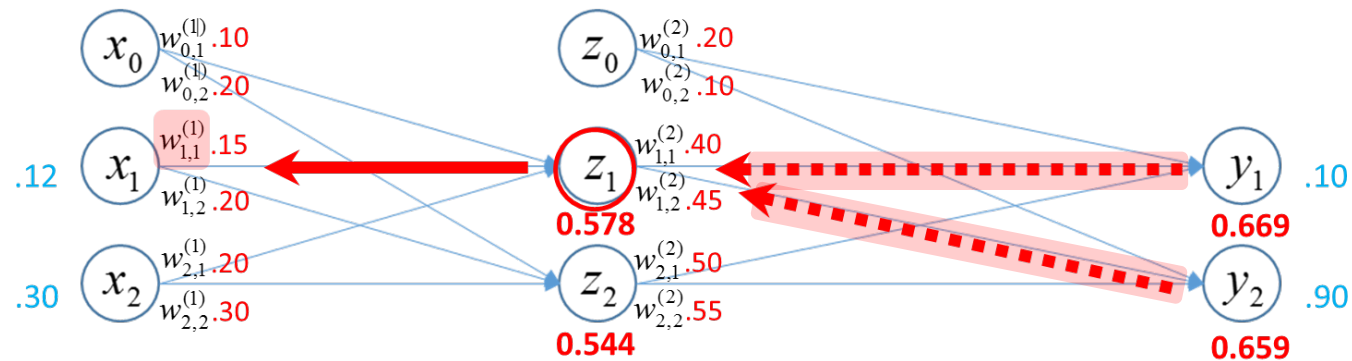
$$\frac{\partial E(w)}{\partial w_{1,1}^{(2)}} = \mathbf{0.0727}$$

Backpropagation algorithm - Backwarding



$$\frac{\partial z_{1(input)}}{\partial w_{1,1}^{(1)}} \times \frac{\partial z_{1(output)}}{\partial z_{1(input)}} \times \frac{\partial E(w)}{\partial z_{1(output)}} = \frac{\partial E(w)}{\partial w_{1,1}^{(1)}}$$

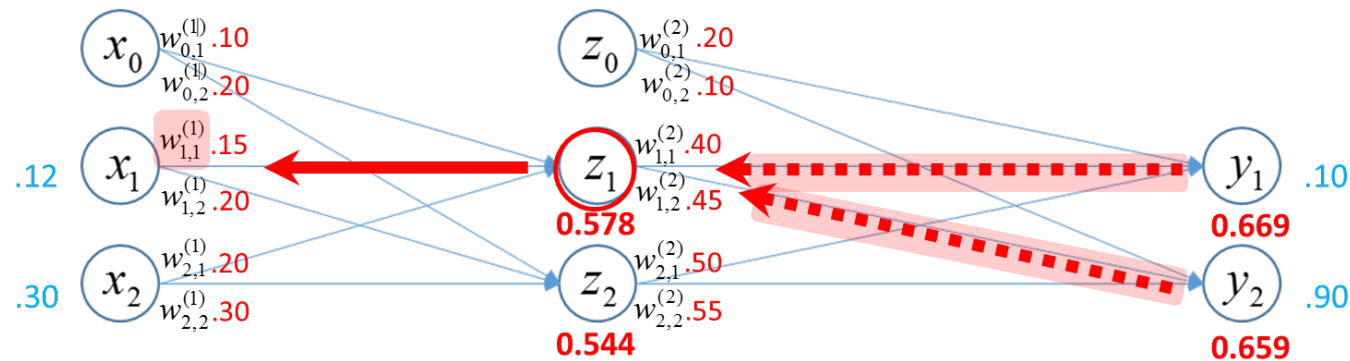
Backpropagation algorithm - Backwarding



$$\frac{\partial z_{1(input)}}{\partial w_{1,1}^{(1)}} \times \frac{\partial z_{1(output)}}{\partial z_{1(input)}} \times \frac{\partial E(w)}{\partial z_{1(output)}} = \frac{\partial E(w)}{\partial w_{1,1}^{(1)}}$$

$$\frac{\partial E(w)}{\partial z_{1(output)}} = \frac{\partial E(w)_{y_1}}{\partial z_{1(output)}} + \frac{\partial E(w)_{y_2}}{\partial z_{1(output)}}$$

Backpropagation algorithm - Backwarding



$$y_{2(input)} = w_{1,2}^{(2)} z_1 + w_{2,2}^{(2)} z_2 + w_{0,2}^{(2)}$$

$$\frac{\partial y_{2(input)}}{\partial z_{1(output)}} = w_{1,2}^{(2)} = 0.45$$

$$\frac{\partial z_{1(input)}}{\partial w_{1,1}^{(1)}} \times \frac{\partial z_{1(ouput)}}{\partial z_{1(input)}} \times \frac{\partial E(w)}{\partial z_{1(ouput)}} = \frac{\partial E(w)}{\partial w_{1,1}^{(1)}}$$

0.0259

$$\frac{\partial E(w)}{\partial z_{1(ouput)}} = \frac{\partial E(w)_{y_1}}{\partial z_{1(ouput)}} + \frac{\partial E(w)_{y_2}}{\partial z_{1(ouput)}}$$

$$= 0.0503 - 0.0244 = 0.0259$$

Previously calculated

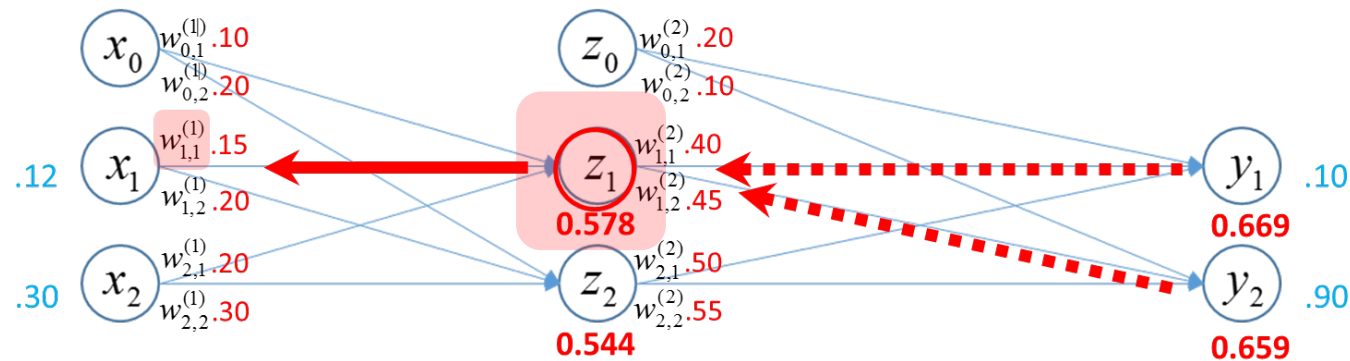
$$\frac{\partial E(w)_{y_1}}{\partial z_{1(ouput)}} = \frac{\partial E(w)_{y_1}}{\partial y_{1(output)}} \times \frac{\partial y_{1(output)}}{\partial y_{1(input)}} \times \frac{\partial y_{1(input)}}{\partial z_{1(ouput)}}$$

0.0503 0.569 0.221 0.4

$$\frac{\partial E(w)_{y_2}}{\partial z_{1(ouput)}} = \frac{\partial E(w)_{y_2}}{\partial y_{2(output)}} \times \frac{\partial y_{2(output)}}{\partial y_{2(input)}} \times \frac{\partial y_{2(input)}}{\partial z_{2(ouput)}}$$

-0.0244 -0.241 0.225 0.45

Backpropagation algorithm - Backwarding



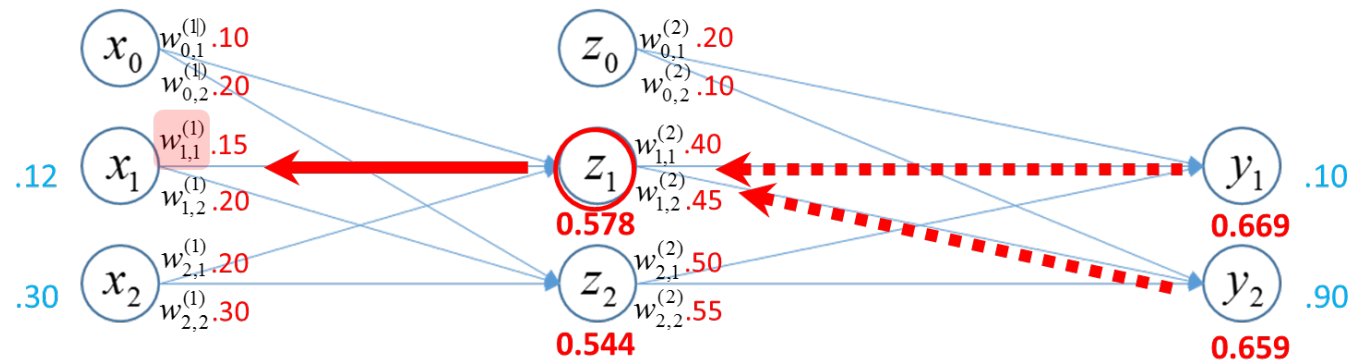
$$\frac{\partial z_{1(input)}}{\partial w_{1,1}^{(1)}} \times \frac{\partial z_{1(output)}}{\partial z_{1(input)}} \times \frac{\partial E(w)}{\partial z_{1(output)}} = \frac{\partial E(w)}{\partial w_{1,1}^{(1)}}$$

0.2439 **0.0259**

$$z_{1(output)} = \frac{1}{1 + e^{-z_{1(input)}}}$$

$$\begin{aligned} \frac{\partial z_{1(output)}}{\partial z_{1(input)}} &= \sigma(z_{1(input)}) (1 - \sigma(z_{1(input)})) \\ &= z_{1(output)} (1 - z_{1(output)}) \\ &= 0.578 \times (1 - 0.578) = 0.2439 \end{aligned}$$

Backpropagation algorithm - Backwarding

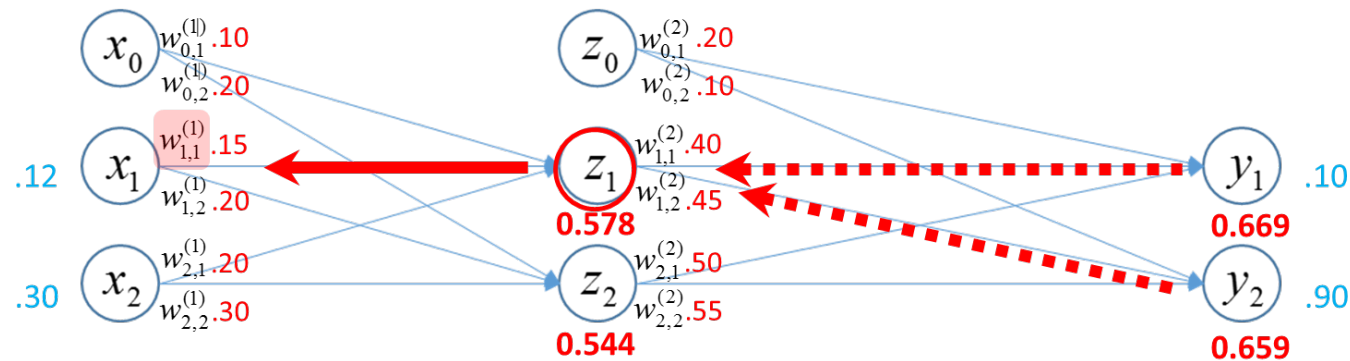


$$\frac{\frac{\partial z_{1(input)}}{\partial w_{1,1}^{(1)}}}{0.12} \times \frac{\frac{\partial z_{1(output)}}{\partial z_{1(input)}}}{0.2439} \times \frac{\frac{\partial E(w)}{\partial z_{1(output)}}}{0.0259} = \frac{\partial E(w)}{\partial w_{1,1}^{(1)}}$$

$$z_{1(input)} = w_{1,1}^{(1)} x_1 + w_{2,1}^{(1)} x_2 + w_{0,1}^{(1)}$$

$$\frac{\partial z_{1(input)}}{\partial w_{1,1}^{(1)}} = x_1 = 0.12$$

Backpropagation algorithm - Backwarding



$$\frac{\partial z_{1(input)}}{\partial w_{1,1}^{(1)}} \times \frac{\partial z_{1(output)}}{\partial z_{1(input)}} \times \frac{\partial E(w)}{\partial z_{1(output)}} = \frac{\partial E(w)}{\partial w_{1,1}^{(1)}}$$

0.12 0.2439 0.0259

$$w_{1,1}^{(1)*} = w_{1,1}^{(1)} + \eta \frac{\partial E(w)}{\partial w_{1,1}^{(1)}}$$

$$= 0.15 + 0.5 \times 0.00075$$

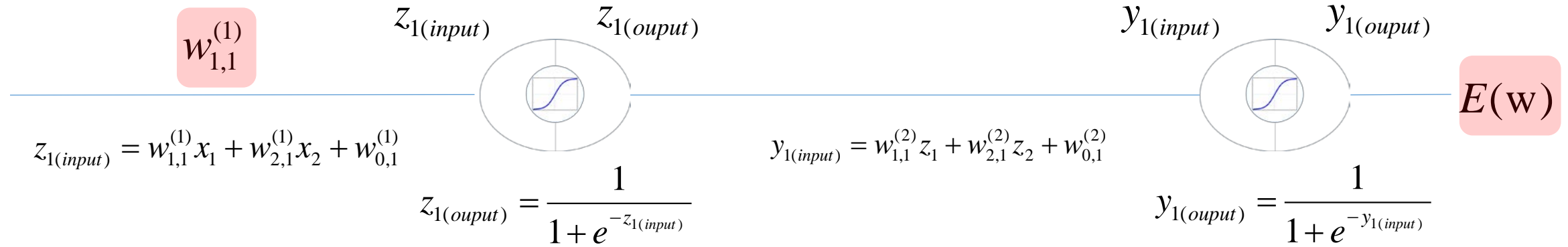
$$= \mathbf{0.1504}$$

$$\frac{\partial E(w)}{\partial w_{1,1}^{(1)}} = \mathbf{0.00075}$$

Same procedures are applied for

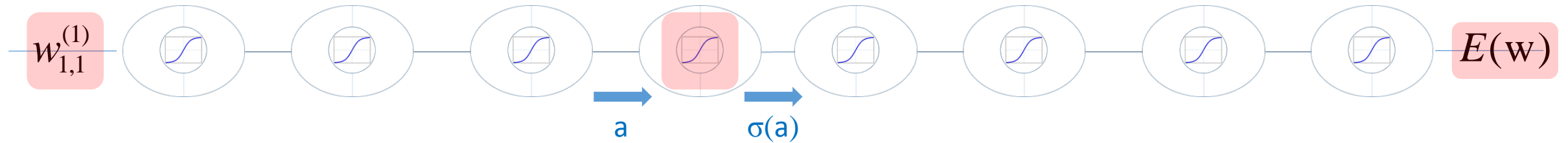
$$w_{1,2}^{(1)}, w_{2,1}^{(1)}, w_{2,2}^{(1)}$$

Backpropagation algorithm - Backwarding

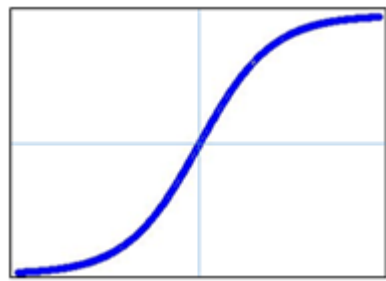


$$\frac{\partial z_{1(input)}}{\partial w_{1,1}^{(1)}} \times \frac{\partial z_{1(output)}}{\partial z_{1(input)}} \times \frac{\partial y_{1(output)}}{\partial y_{1(input)}} \times \frac{\partial E(w)}{\partial y_{1(output)}} = \frac{\partial E(w)}{\partial w_{1,1}^{(1)}}$$

Vanishing gradient problem



$$\frac{\partial E(w)}{\partial w_{1,1}^{(1)}} = \frac{\partial E(w)}{\partial z_{1(ouput)}} \times \dots \times \frac{\partial \sigma(a)}{\partial a} \times \dots \times \frac{\partial z_{1(input)}}{\partial w_{1,1}^{(1)}}$$



$$\frac{\partial \sigma(a)}{\partial a} \approx 0$$

$$w_{new} = w_{old} + \eta \frac{\partial E(w)}{\partial w_{1,1}^{(1)}}$$

$$\frac{\partial E(w)}{\partial w_{1,1}^{(1)}} \approx 0$$

Gradient
vanishing

Backup Slides