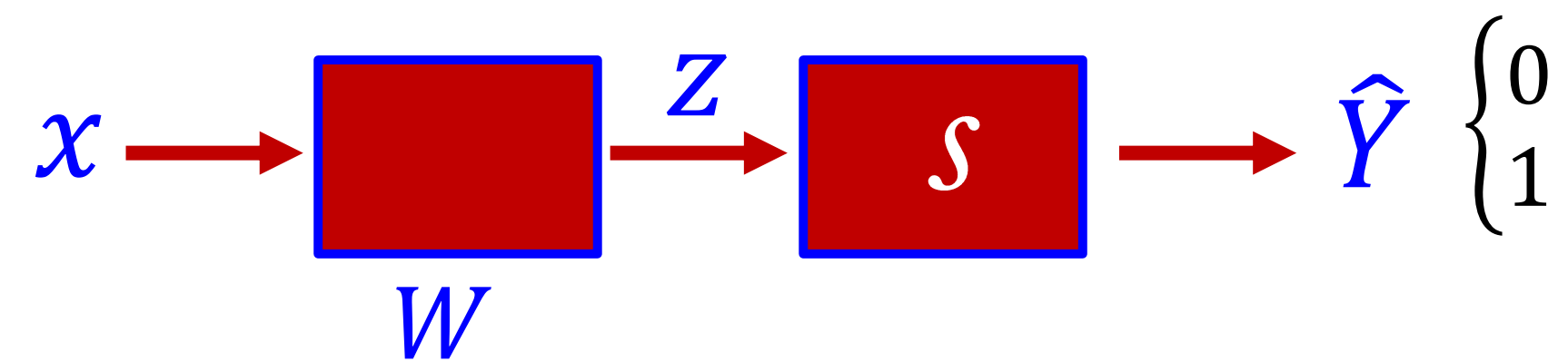


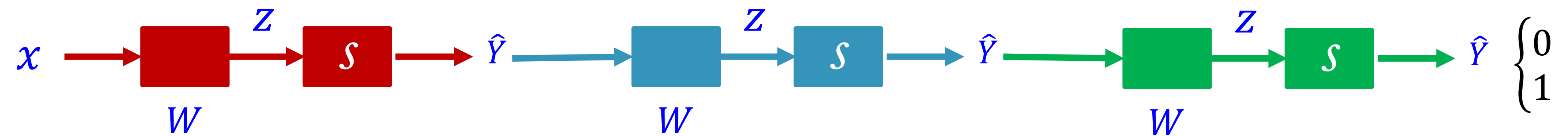
Lecture 09

# Neural Nets(NN) for XOR, Back propagation

# One logistic regression unit **cannot** separate XOR

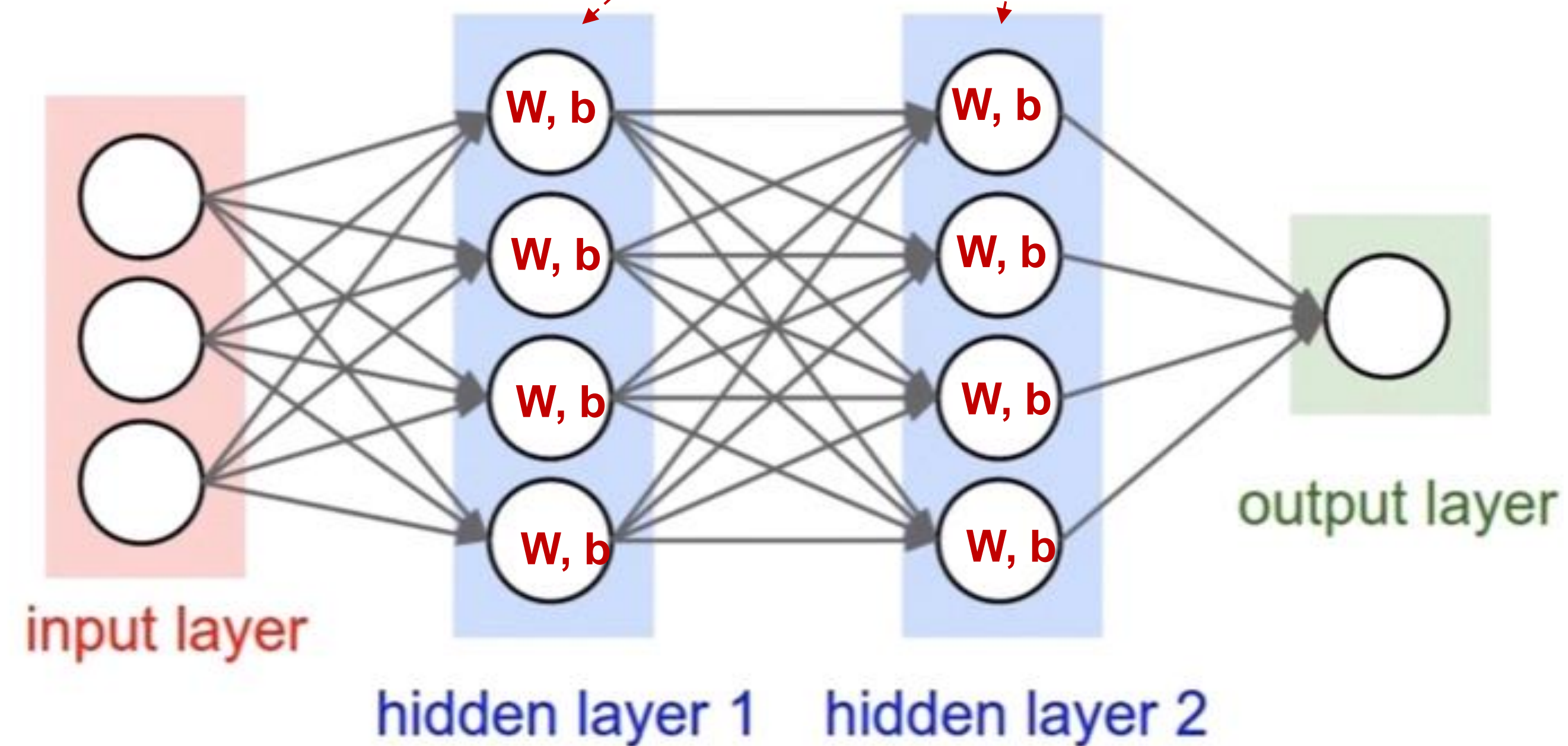


# Multiple logistic regression units



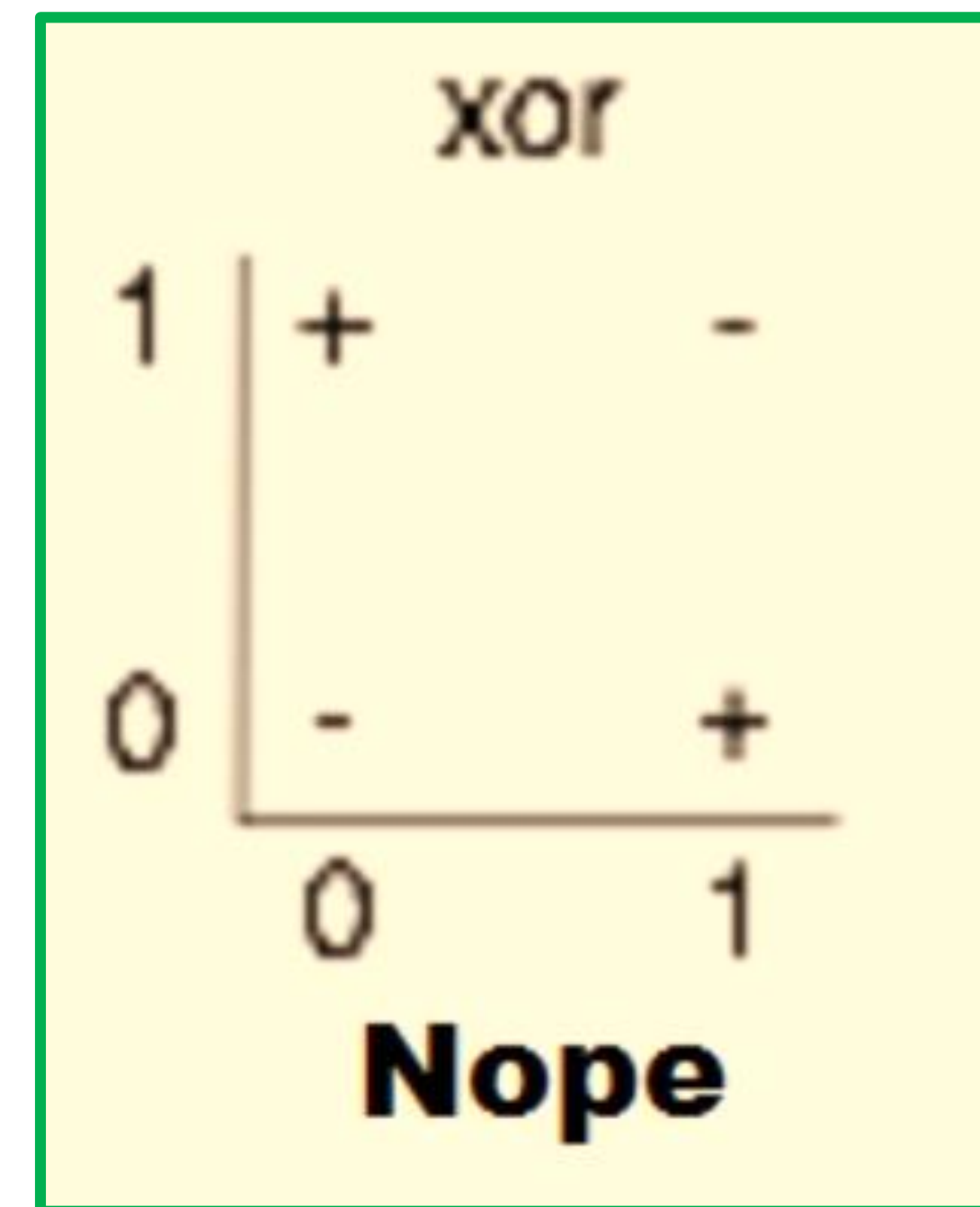
Neural Network (NN) : “No one on earth had found a viable way to train”

[Marvin Minsky]



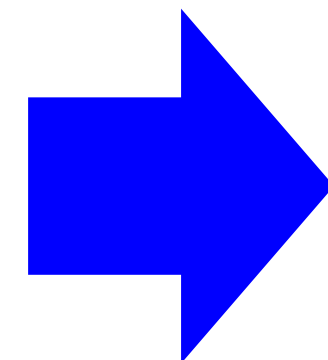
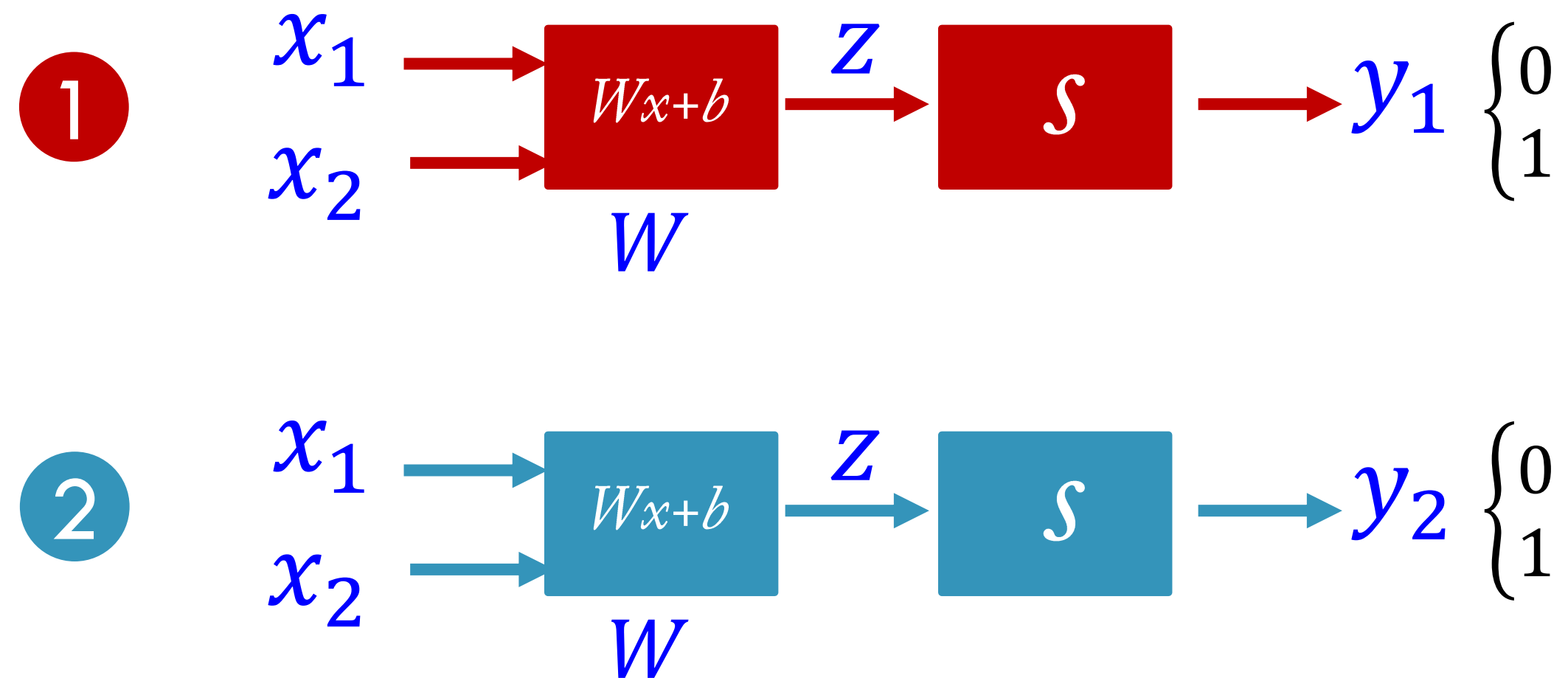
# XOR using NN

<b>X<sub>1</sub></b>	<b>X<sub>2</sub></b>	<b>XOR</b>
<b>0</b>	<b>0</b>	<b>0 (-)</b>
<b>0</b>	<b>1</b>	<b>1 (+)</b>
<b>1</b>	<b>0</b>	<b>1 (+)</b>
<b>1</b>	<b>1</b>	<b>0 (-)</b>

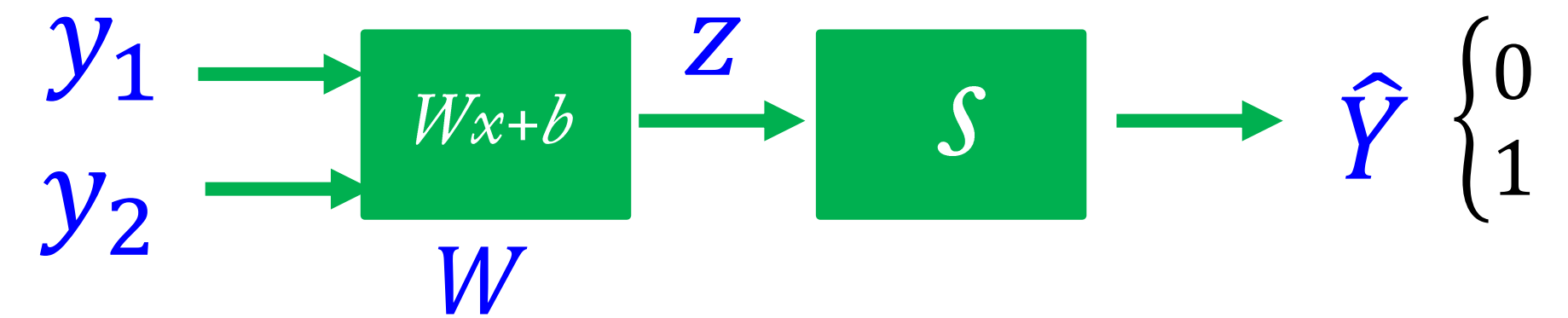


# Neural Net

$$H(x) = Wx + b$$



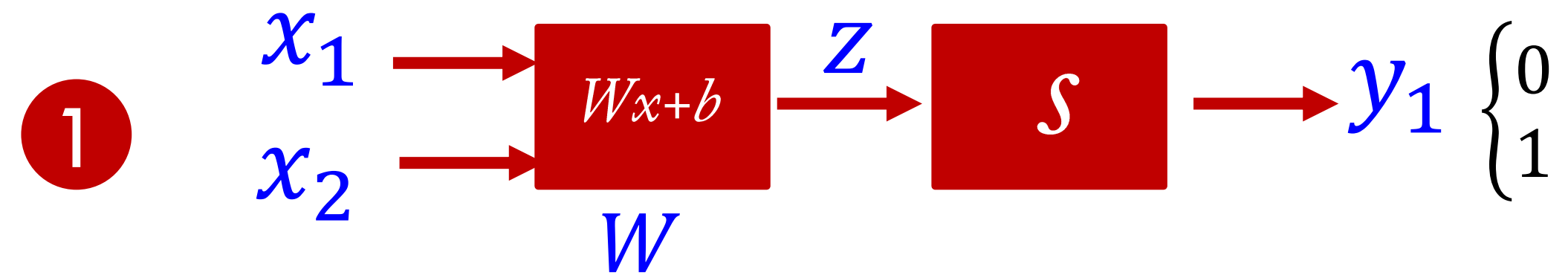
3



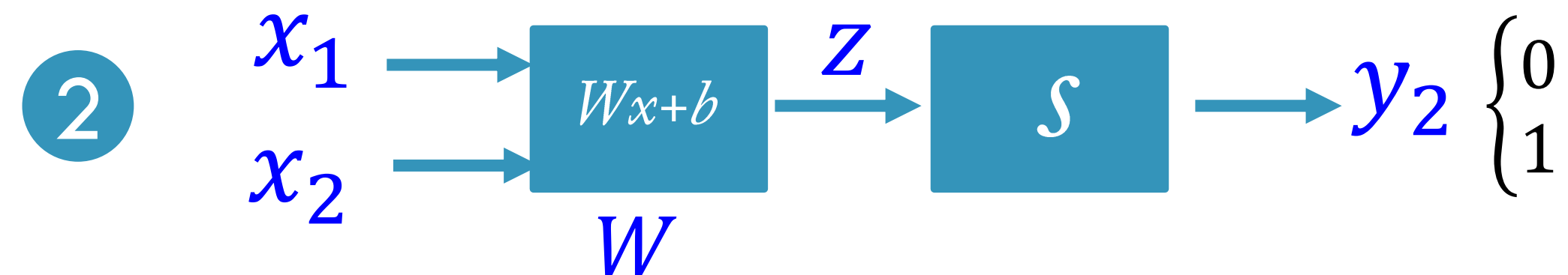
# Neural Net

$$H(x) = Wx + b$$

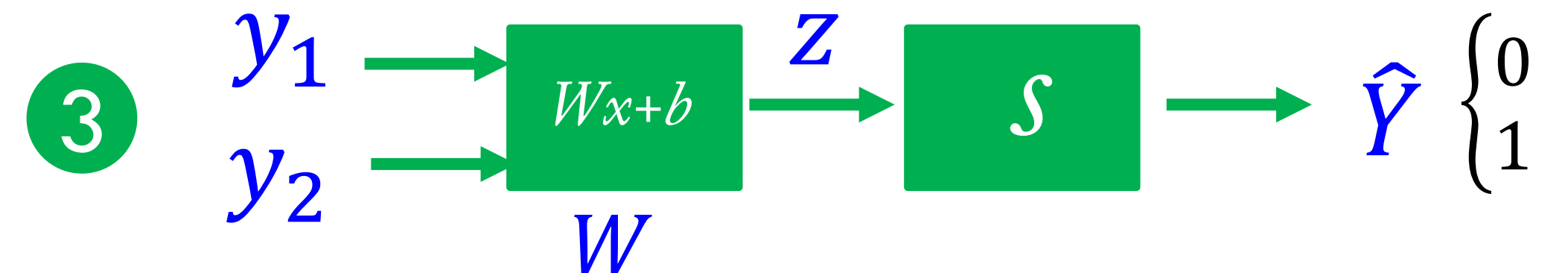
if)  $w = \begin{bmatrix} 5 \\ 5 \end{bmatrix}, b = -8$



if)  $w = \begin{bmatrix} -7 \\ -7 \end{bmatrix}, b = 3$

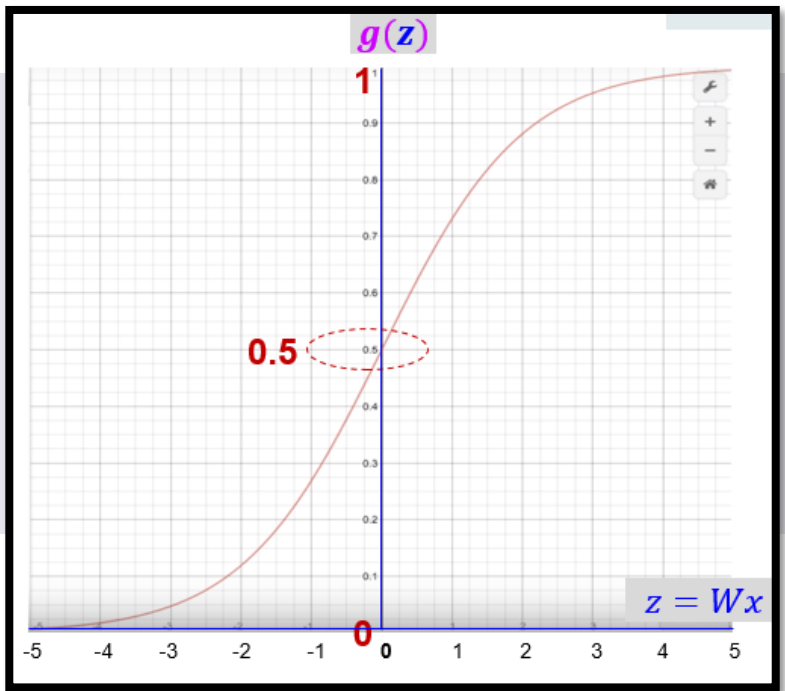


if)  $w = \begin{bmatrix} -1 & 1 \\ -1 & 1 \end{bmatrix}, b = 6$



# Neural Net

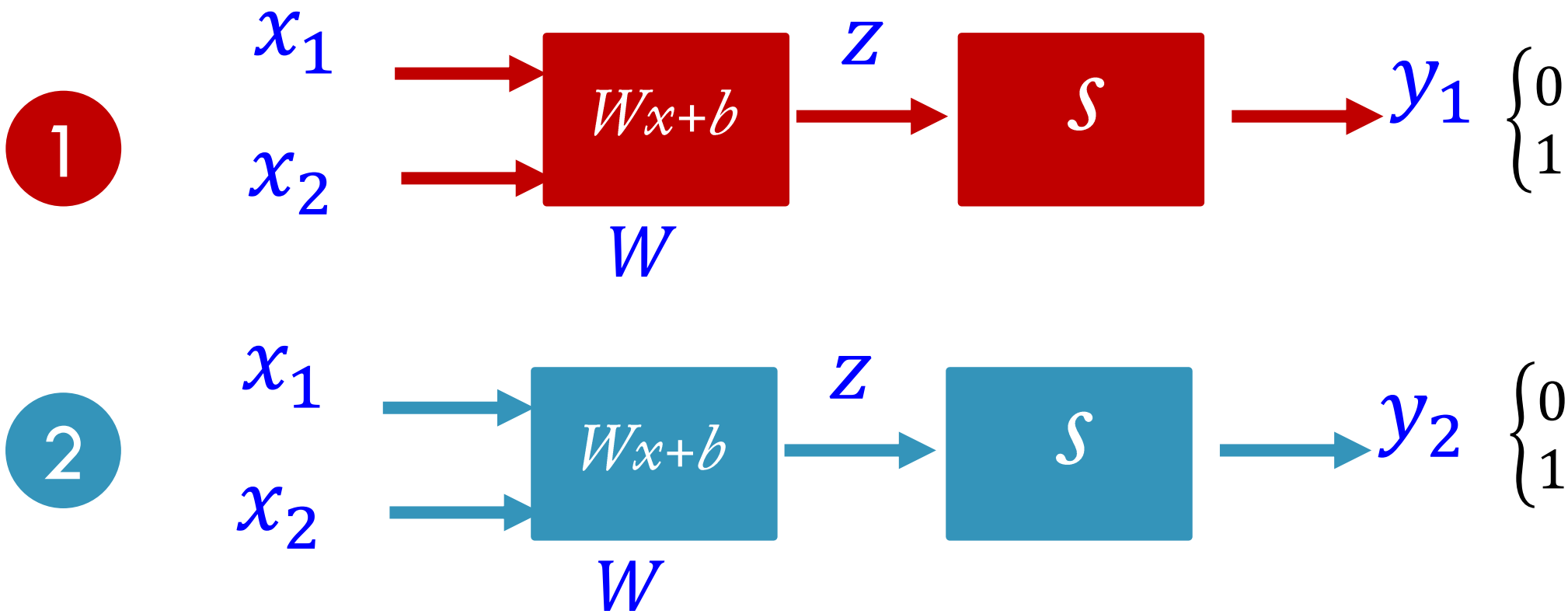
$x_1$	$x_2$	XOR
0	0	0 (-)
0	1	1 (+)
1	0	1 (+)
1	1	0 (-)



$$H(x) = Wx + b$$

$W=(5, 5), b=-8$  일때  
 $(0*5 + 0*5) + (-8) = -8 \implies \text{sigmoid}(-8) = 0$   
 $(0*5 + 1*5) + (-8) = -3 \implies \text{sigmoid}(-3) = 0$   
 $(1*5 + 0*5) + (-8) = -3 \implies \text{sigmoid}(-3) = 0$   
 $(1*5 + 1*5) + (-8) = 2 \implies \text{sigmoid}(2) = 1$

$$w = \begin{bmatrix} 5 \\ 5 \end{bmatrix}, b = -8$$

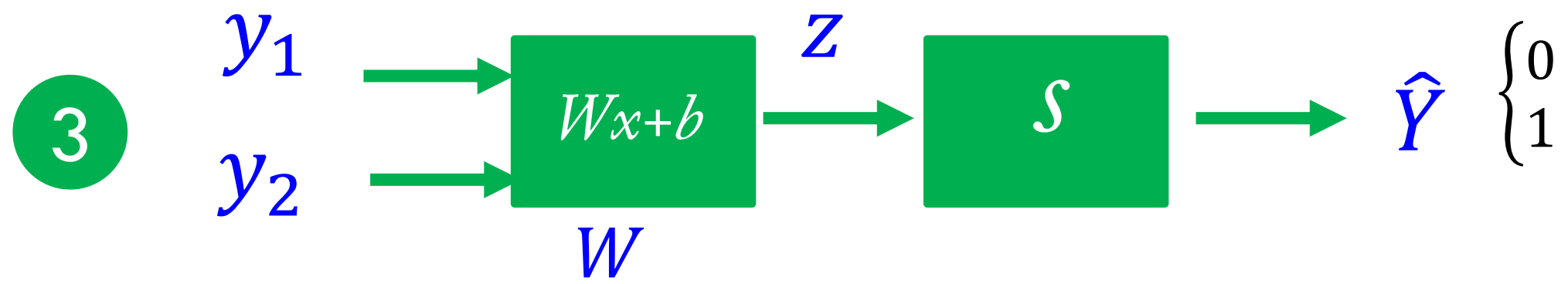


$W=(-7, -7), b=3$  일때,  
 $(0*-7 + 0*-7) + 3 = 3 \implies \text{sigmoid}(3) = 1$   
 $(0*-7 + 1*-7) + 3 = -4 \implies \text{sigmoid}(-4) = 0$   
 $(1*-7 + 0*-7) + 3 = -4 \implies \text{sigmoid}(-4) = 0$   
 $(1*-7 + 1*-7) + 3 = -11 \implies \text{sigmoid}(-11) = 0$

$$w = \begin{bmatrix} -7 \\ -7 \end{bmatrix}, b = 3$$

$$w = \begin{bmatrix} -11 \\ -11 \end{bmatrix}, b = 6$$

$W=(-11, -11), b=6$  일때  
 $(0*-11 + 1*-11) + 6 = -5 \implies \text{sigmoid}(-5) = 0$   
 $(0*-11 + 0*-11) + 6 = 6 \implies \text{sigmoid}(6) = 1$   
 $(0*-11 + 0*-11) + 6 = 6 \implies \text{sigmoid}(6) = 1$   
 $(1*-11 + 0*-11) + 6 = -5 \implies \text{sigmoid}(-5) = 0$

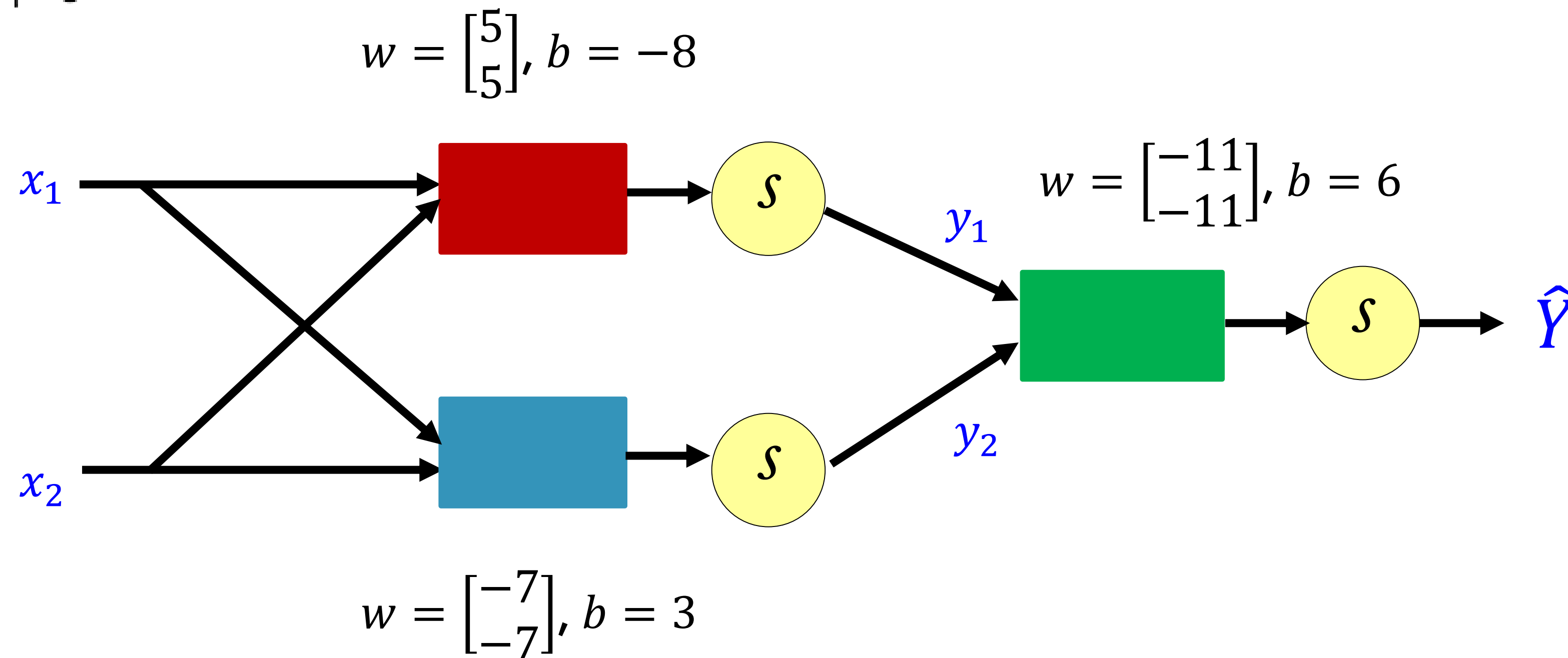


$x_1$	$x_2$	$y_1$	$y_2$	$\hat{y}$	XOR
0	0	0	1	0	0 (-)
0	1	0	0	1	1 (+)
1	0	0	0	1	1 (+)
1	1	1	0	0	0 (-)



# Forward propagation

$$H(x) = Wx + b$$



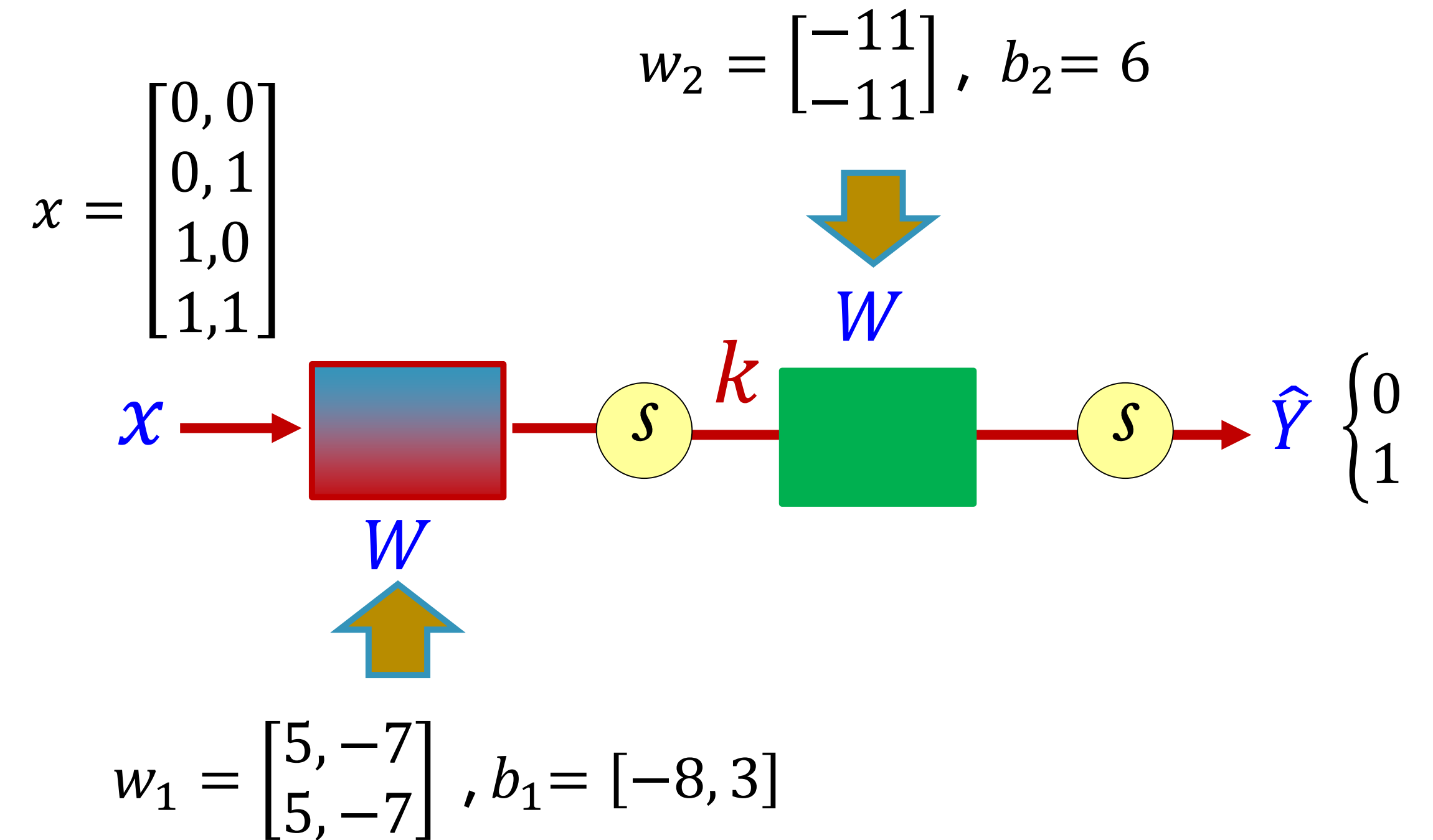
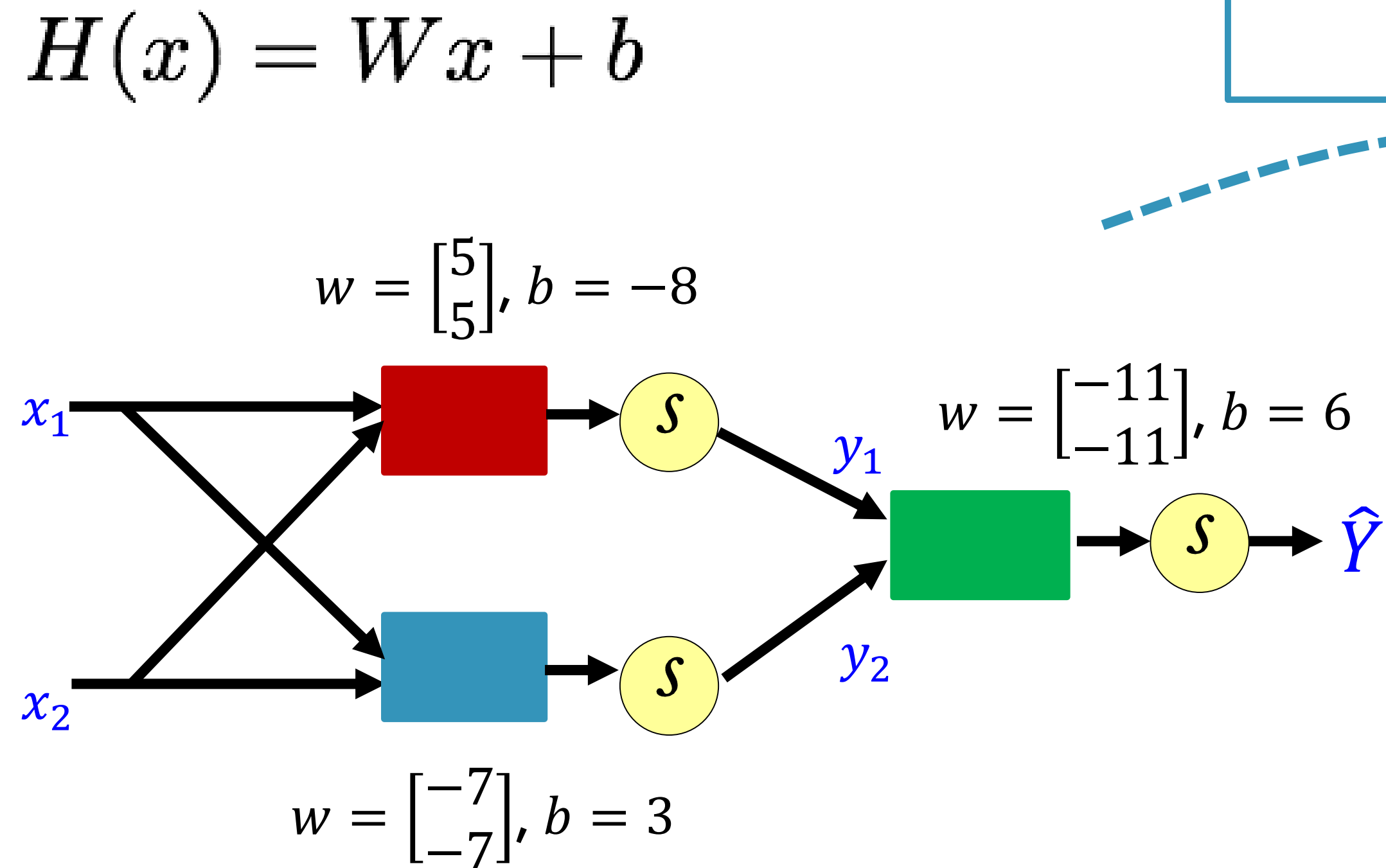
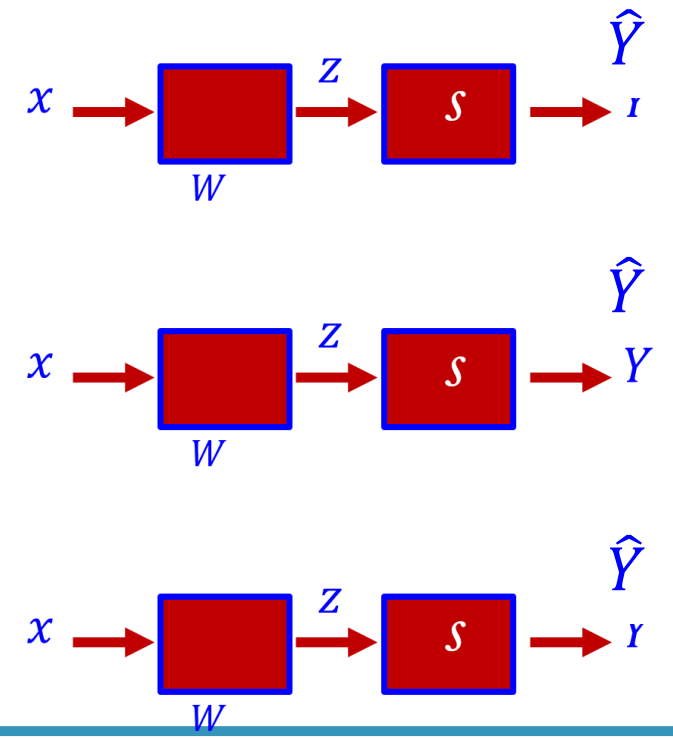
*Can you find another  $W$  and  $b$  for the XOR?*

# Neural Network (NN)

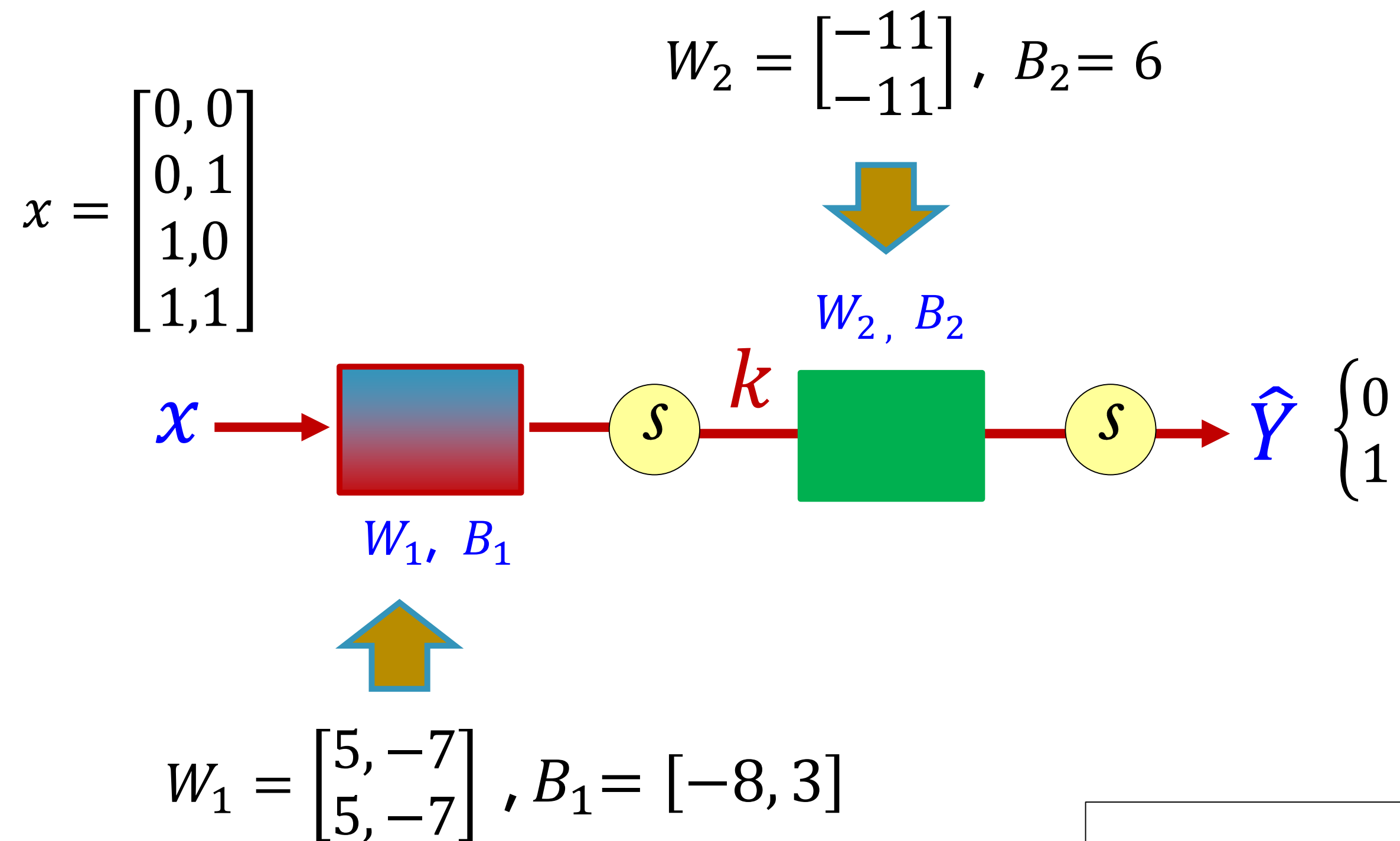
## Multinomial classification

$$\begin{bmatrix} w_{A1} & w_{A2} & w_{A3} \\ w_{B1} & w_{B2} & w_{B3} \\ w_{C1} & w_{C2} & w_{C3} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} w_{A1}x_1 + w_{A2}x_2 + w_{A3}x_3 \\ w_{B1}x_1 + w_{B2}x_2 + w_{B3}x_3 \\ w_{C1}x_1 + w_{C2}x_2 + w_{C3}x_3 \end{bmatrix} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \hat{y}_3 \\ C \end{bmatrix}$$

$H_A(x)$   
 $H_B(x)$   
 $H_C(x)$



# Neural Network (NN)



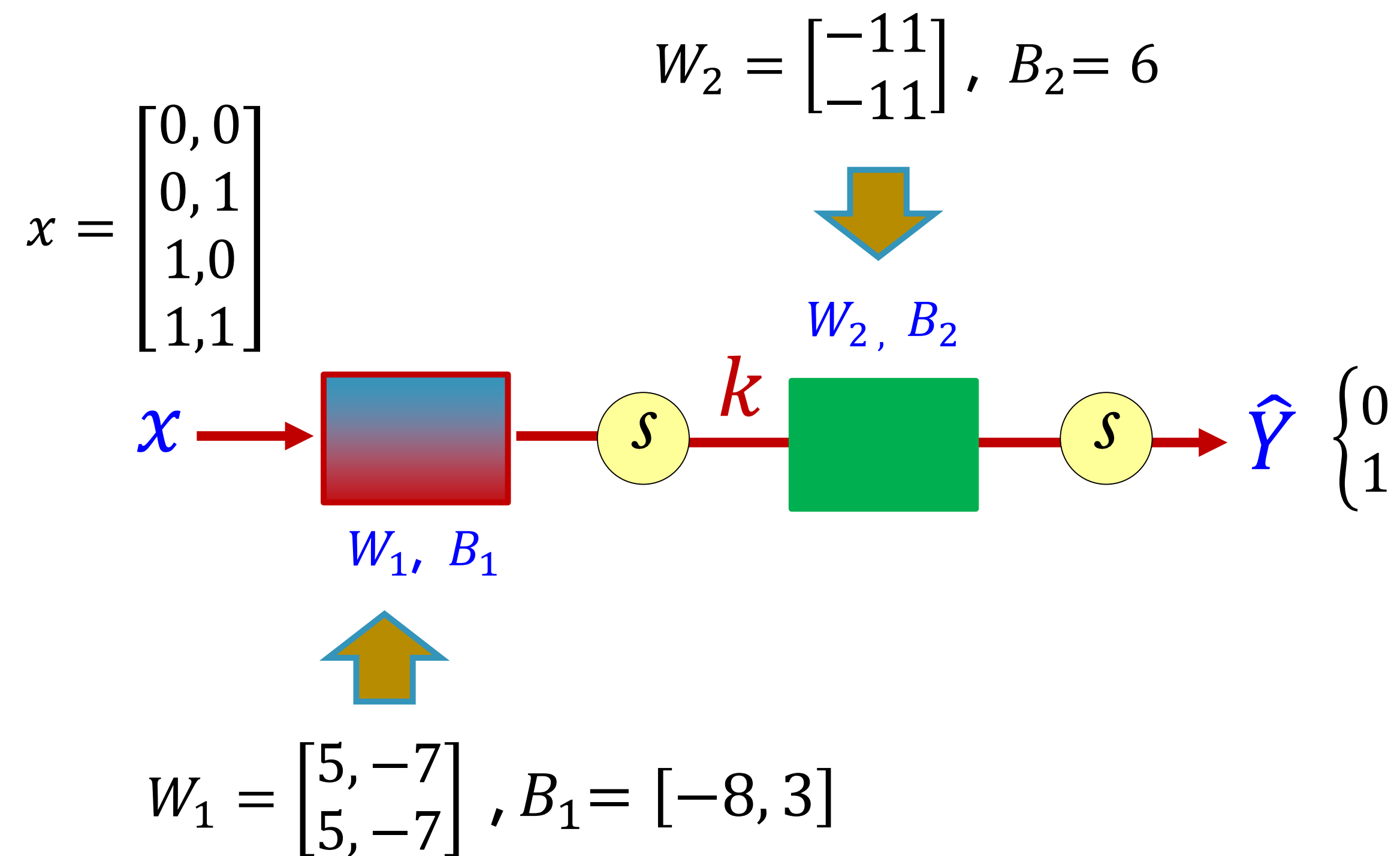
$$H(x) = Wx + b$$

$$k(x) = \text{sigmoid}(x \cdot W_1 + B_1)$$

$$\hat{Y} = H(x) = \text{Sigmoid}(k(x) \cdot W_2 + B_2)$$

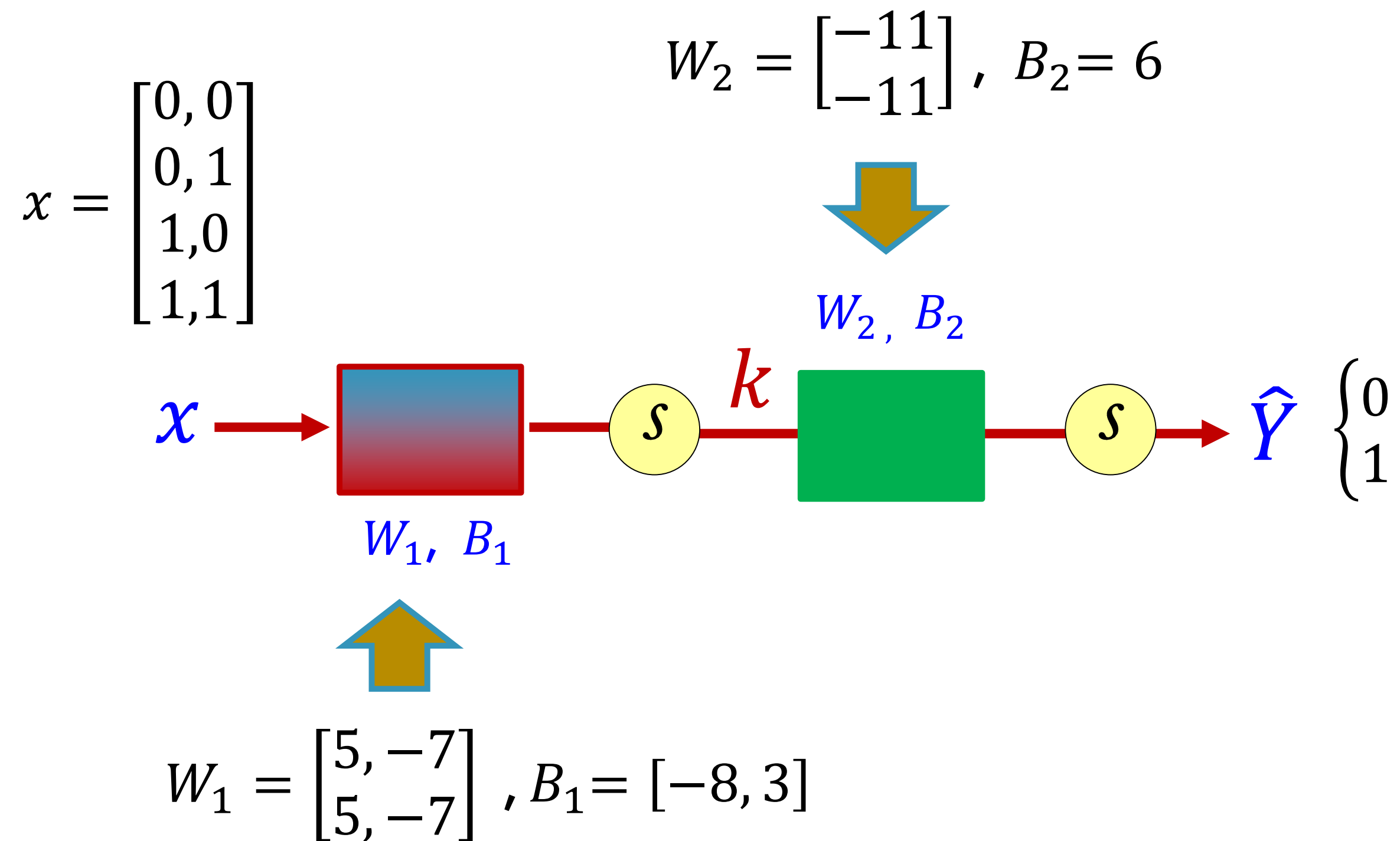
```
# NN
K = tf.sigmoid(tf.matmul(X, W1) + b1)
hypothesis = tf.sigmoid(tf.matmul(K, W2) + b2)
```

# Neural Network (NN)



*How can we learn  $W_1, W_2, B_1, b_2$  from training data?*

# Neural Network (NN)



$$H(x) = Wx + b$$

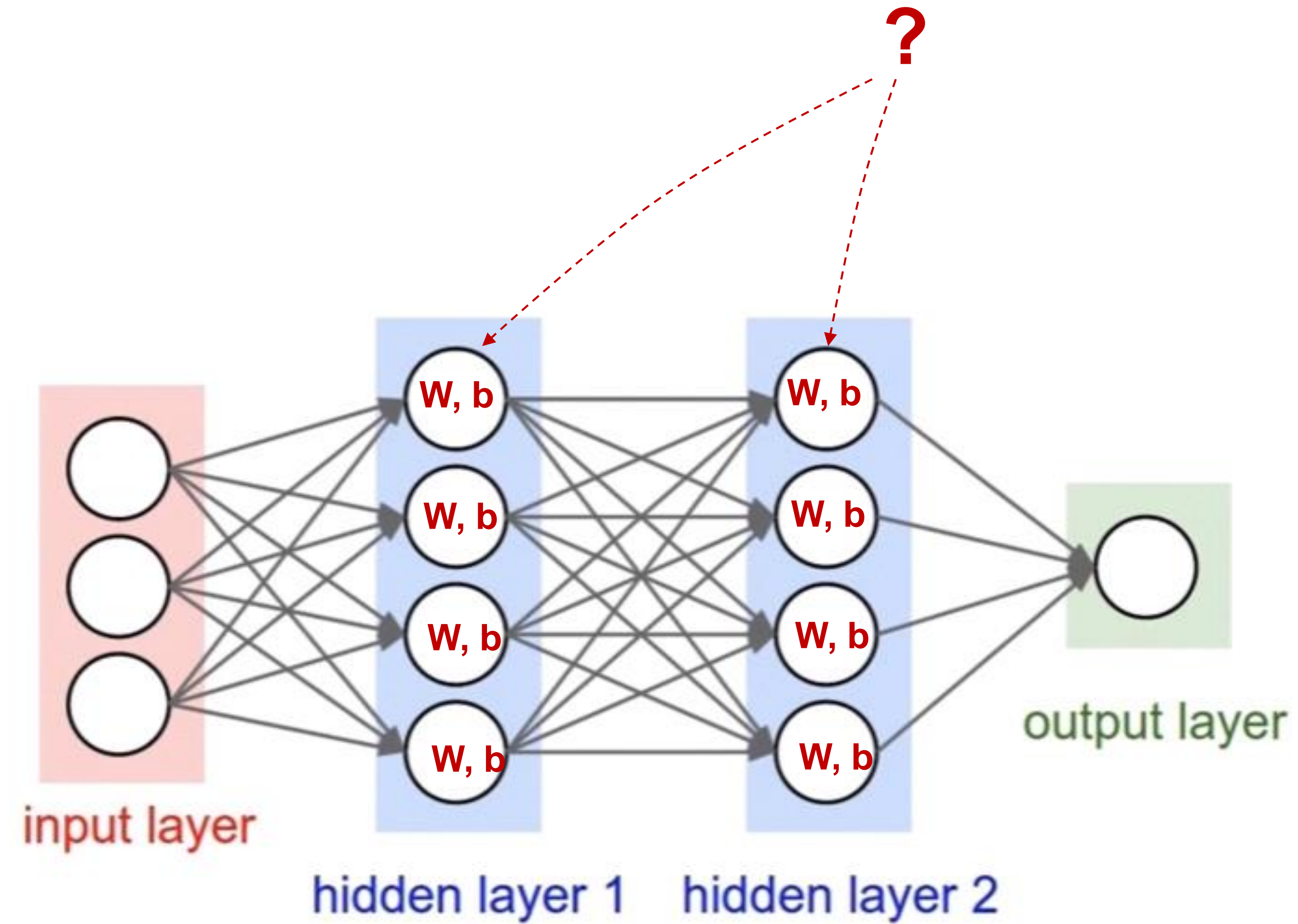
$$k(x) = \text{sigmoid}(x \cdot W_1 + B_1)$$

$$\hat{Y} = H(x) = \text{Sigmoid}(k(x) \cdot W_2 + B_2)$$

How can we learn  $W_1, W_2, B_1, b_2$  from training data?

**Gradient Descent**

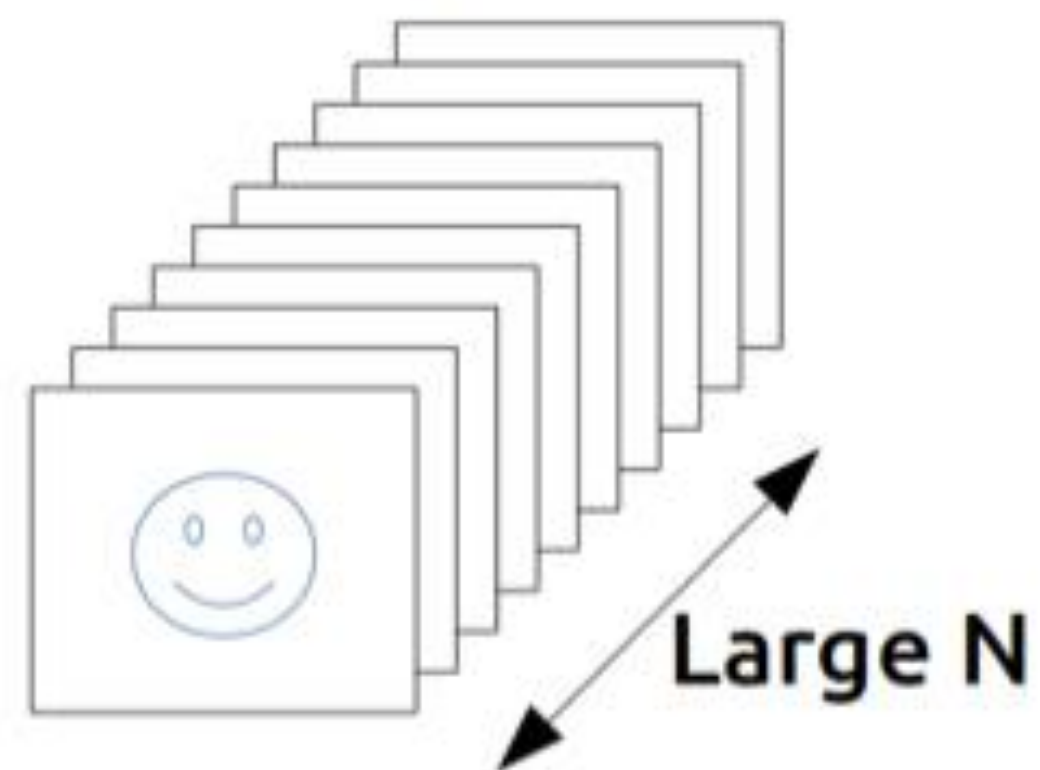
# Derivation



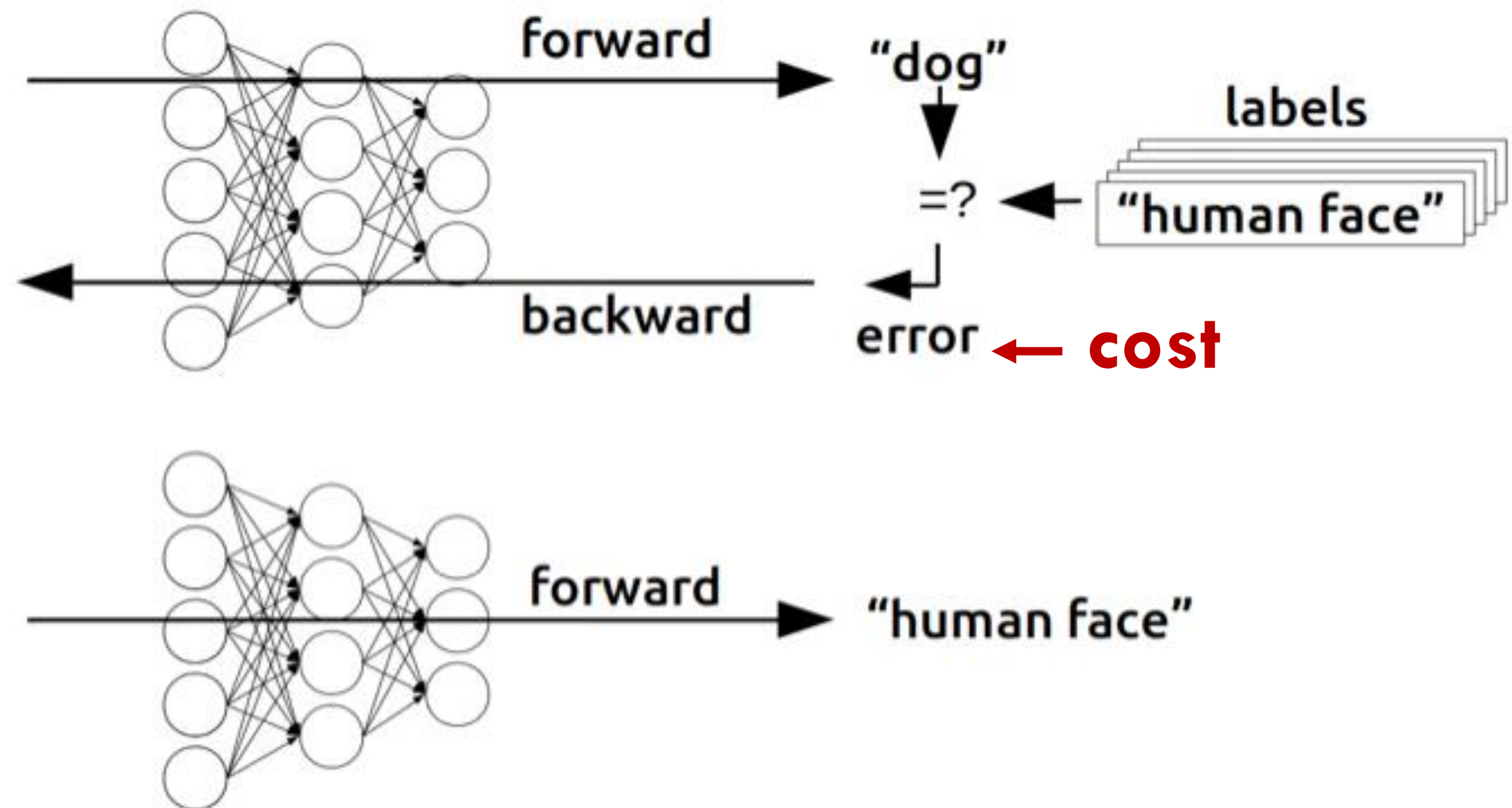
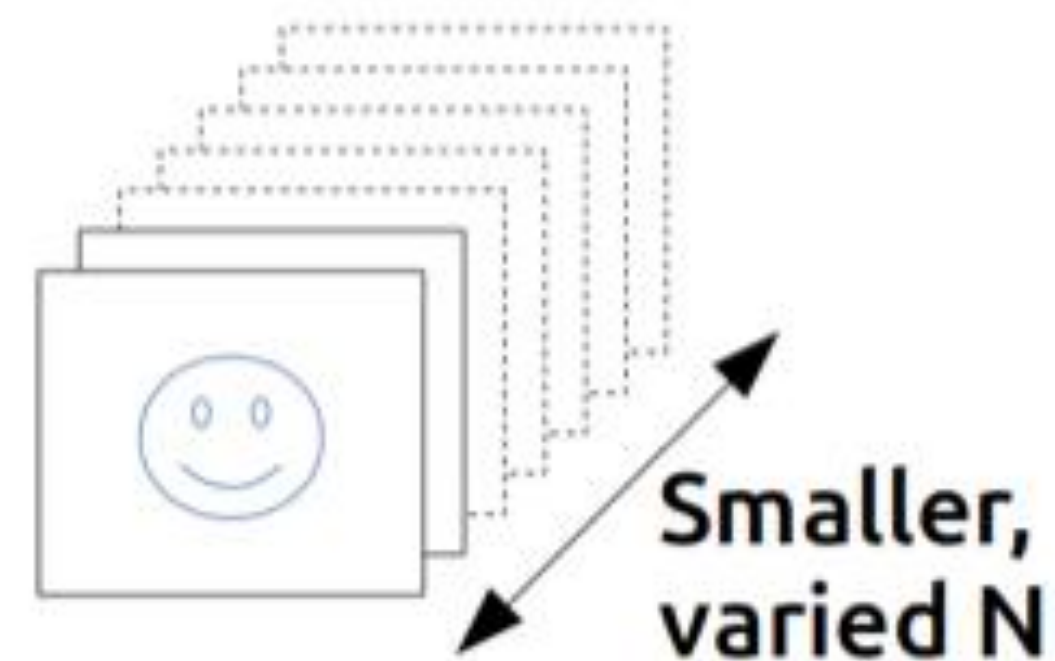


# Backpropagation (1974, 1982 by Paul Werbos, 1986 by Hinton)

## Training



## Inference



# Basic Derivative

$$\frac{d}{dx}f(x) = \lim_{\Delta x \rightarrow 0} \frac{f(x + \Delta x) - f(x)}{\Delta x} = \frac{df}{dx}$$

if )  $\Delta x = 0.01$

**C**  $f(x) = 3$   
 $x : 1 \rightarrow 3$   
 $x : 2 \rightarrow 3$

$$\rightarrow \frac{f(x+0.01) - f(x)}{0.01} = \frac{3-3}{0.01} = \frac{0}{0.01} = 0$$

**V**  $f(x) = x$   
 $x : 1 \rightarrow f(x) : 1$   
 $x : 2 \rightarrow f(x) : 2$

$$\rightarrow \frac{f(x+0.01) - f(x)}{0.01} = \frac{x+0.01-x}{0.01} = \frac{0.01}{0.01} = 1$$

**C** **V**  $f(x) = 2x$   
 $x : 1 \rightarrow f(x) : 2$

$$\rightarrow \frac{f(x+0.01) - f(x)}{0.01} = \frac{2(x+0.01) - 2x}{0.01} = \frac{2x + (2*0.01) - 2x}{0.01} = \frac{2*0.01}{0.01} = 2$$

$f(x) = 3x \rightarrow f(x) : 3$   
 $f(x) = 4x \rightarrow f(x) : 4$

**V** + **C**  $f(x) = x + 3$

$$\rightarrow \frac{f(x+0.01) - f(x)}{0.01} = \frac{(x+0.01+3) - (x+3)}{0.01} = \frac{0.01}{0.01} = 1$$



# Partial Derivative : consider other variables as constants

**C** **V**  $f(x) = 2x \rightarrow \frac{f(x+0.01)-f(x)}{0.01} = \frac{2(x+0.01)-2x}{0.01} = \frac{2x+(2*0.01)-2x}{0.01} = \frac{2*0.01}{0.01} = 2$   
 $x : 1 \rightarrow f(x) : 2$

$f(x) = 3x \rightarrow f(x) : 3$   
 $f(x) = 4x \rightarrow f(x) : 4$

**V** **C**  $f(x, y) = x \cdot y, \frac{\partial f}{\partial x} \rightarrow \text{Partial Derivative} \rightarrow 1 \cdot y = y$   
 상수

**C** **V**  $f(x, y) = x \cdot y, \frac{\partial f}{\partial y} \rightarrow \text{Partial Derivative} \rightarrow x \cdot 1 = x$   
 상수

# Partial Derivative : consider other variables as constants

**V** + **C**  $f(x) = x + 3 \rightarrow \frac{f(x+0.01)-f(x)}{0.01} = \frac{(\cancel{x}+0.01+\cancel{3})-(\cancel{x}+3)}{0.01} = \frac{0.01}{0.01} = 1$

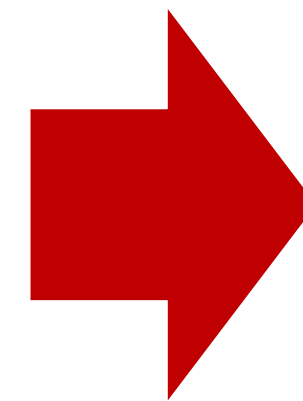
**V** + **C**  $f(x, y) = x + \overset{\text{상수}}{y}, \frac{\partial f}{\partial x} \rightarrow \text{Partial Derivative} \rightarrow 1 + 0 = 1$

**C** + **V**  $f(x, y) = \overset{\text{상수}}{x} + y, \frac{\partial f}{\partial y} \rightarrow \text{Partial Derivative} \rightarrow 0 + 1 = 1$

# Chain rule

$$z = W\underline{x} + b$$

$$H = \text{sigmoid}(z)$$



Chain rule

$$f(g(x))$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial g} \cdot \frac{\partial g}{\partial x}$$

Complex function

$$H(x) = Wx + b$$

$$\mathbf{k}(x) = \text{sigmoid}(x \cdot W_1 + B_1)$$

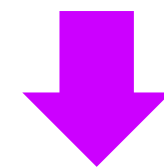
$$\hat{Y} = H(x) = \text{Sigmoid}(\mathbf{k}(x) \cdot W_2 + B_2)$$

```
# NN
K = tf.sigmoid(tf.matmul(X, W1) + b1)
hypothesis = tf.sigmoid(tf.matmul(K, W2) + b2)
```

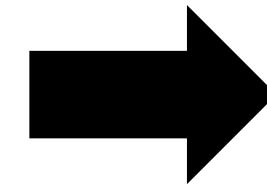
# Back propagation (chain rule)

$$f = wx + b$$

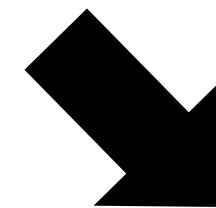
$$g = wx$$



$$f = g + b$$



$$\begin{aligned} & \text{상수} \quad w \cdot x, \frac{\partial g}{\partial w} = 1 \cdot x = x \\ & \text{상수} \quad w \cdot x, \frac{\partial g}{\partial x} = w \cdot 1 = w \end{aligned}$$



$$\begin{aligned} & \text{상수} \quad g + b, \frac{\partial f}{\partial b} \rightarrow 0 + 1 = 1 \\ & \text{상수} \quad g + b, \frac{\partial f}{\partial g} \rightarrow 1 + 0 = 1 \end{aligned}$$

# Back propagation (chain rule)

$$f = wx + b \quad \rightarrow \quad g = wx \quad \rightarrow \quad f = g + b$$

$$\frac{\partial g}{\partial w} = x, \frac{\partial g}{\partial x} = w$$

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

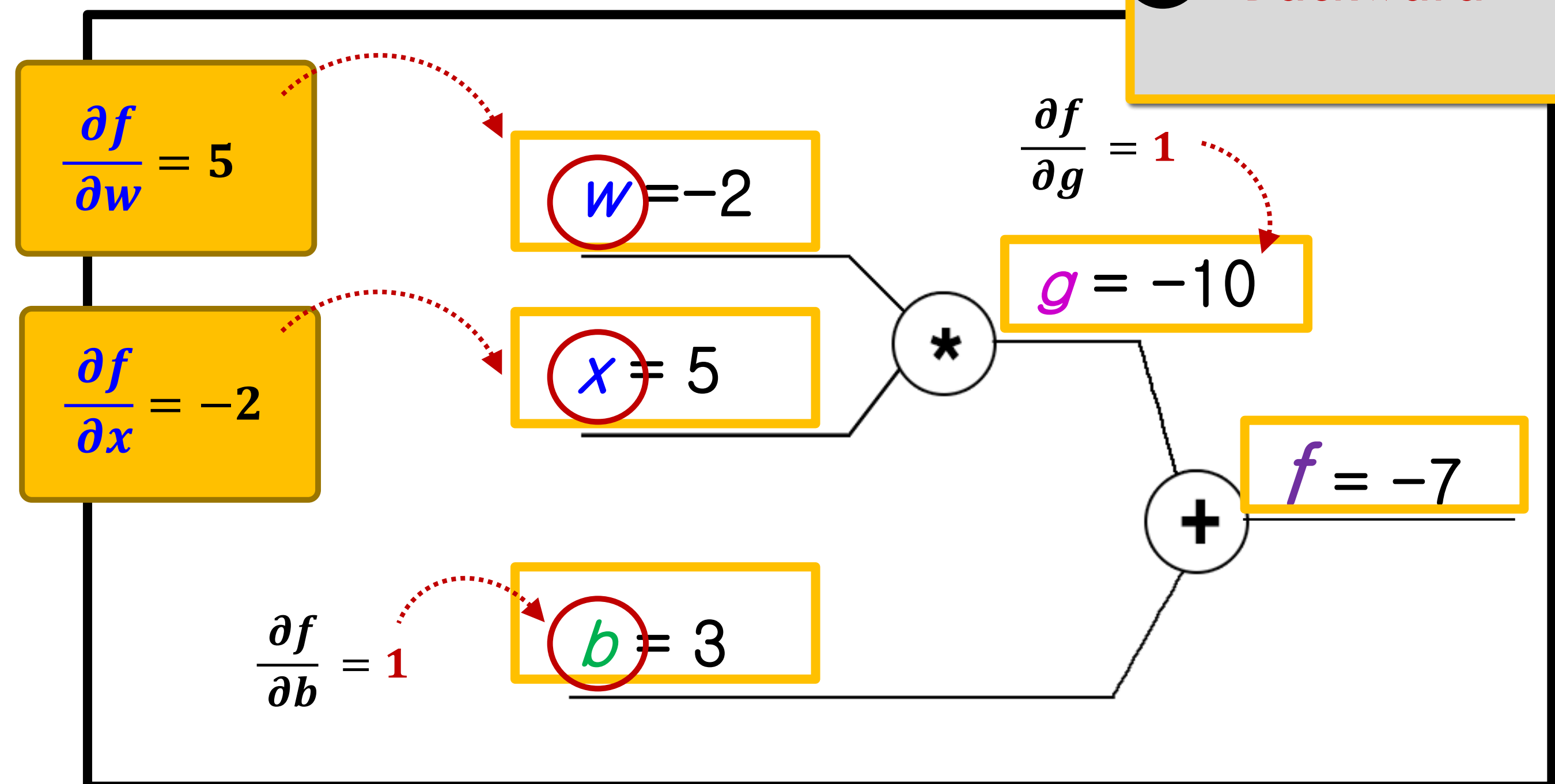
Chain rule

if  
1 Forward  
( $w = -2, x = 5, b = 3$ )  
2 Backward

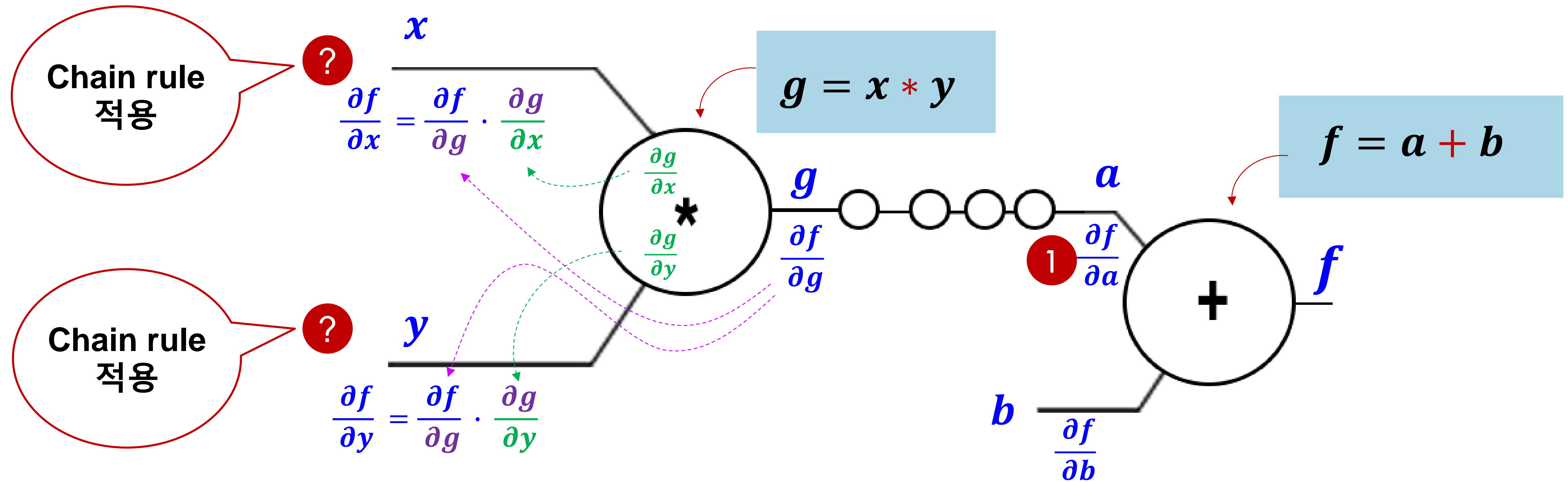
$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial g} \cdot \frac{\partial g}{\partial x}$$

$$\frac{\partial f}{\partial w} = \frac{\partial f}{\partial g} \cdot \frac{\partial g}{\partial w} = 1 \cdot x = 1 \cdot 5 = 5$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial g} \cdot \frac{\partial g}{\partial x} = 1 \cdot w = 1 \cdot -2 = -2$$



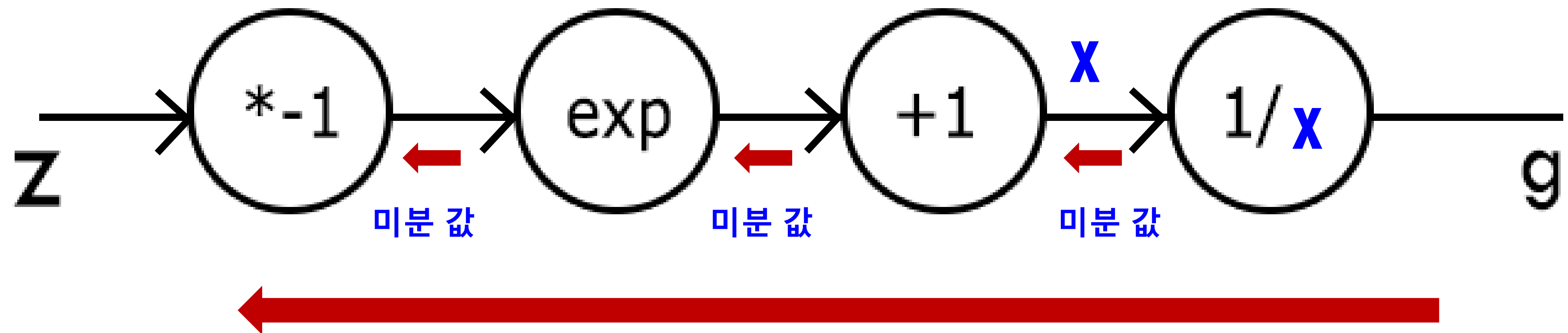
# Back propagation (chain rule)



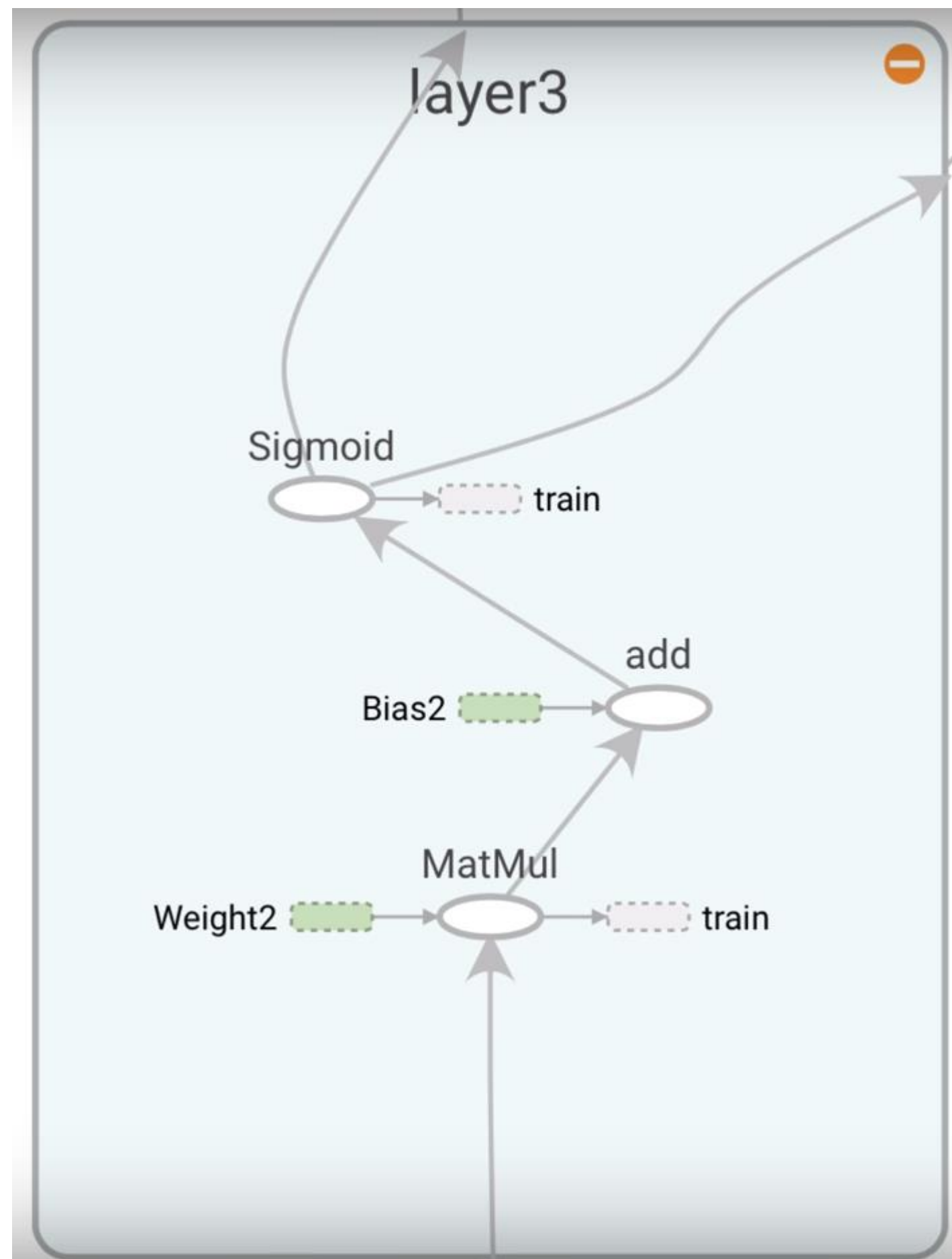
# Back propagation (chain rule) : Sigmoid

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

?  $\frac{\partial g}{\partial z}$



# Back propagation in TensorFlow → TensorBoard



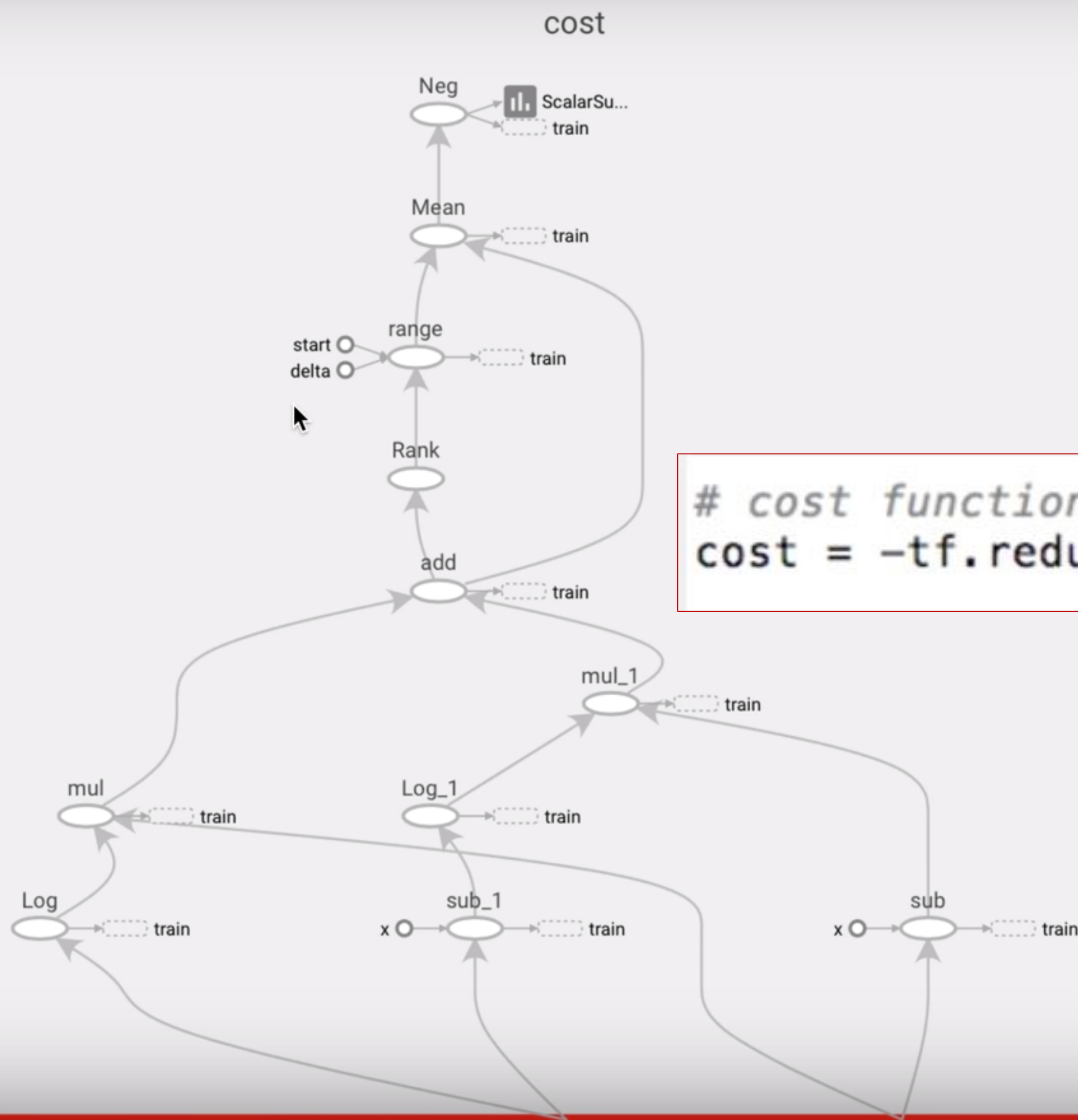
```
hypothesis = tf.sigmoid(tf.matmul(L2, W2) + b2)
```



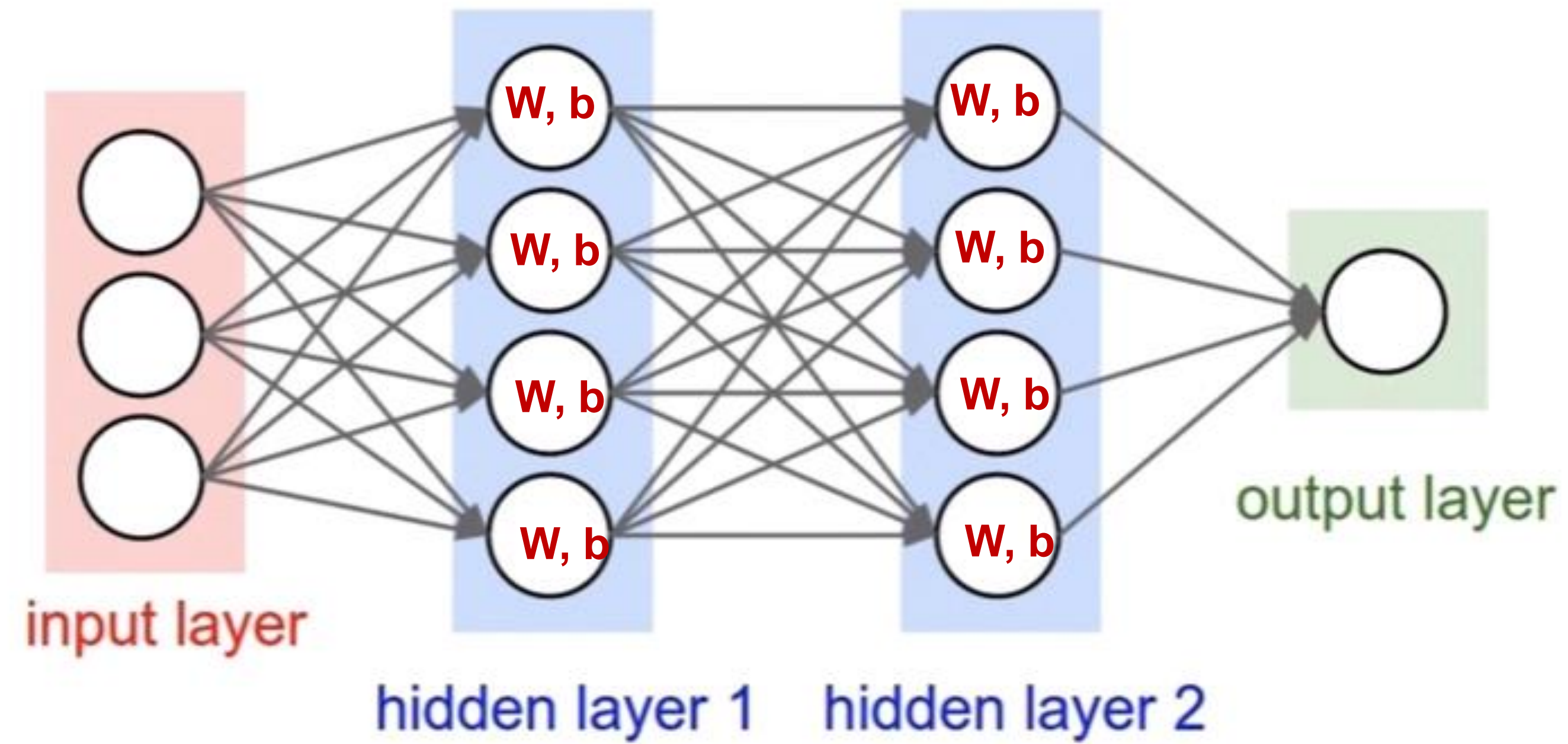
# Back propagation in TensorFlow → TensorBoard

## Back propagation in TensorFlow TensorBoard

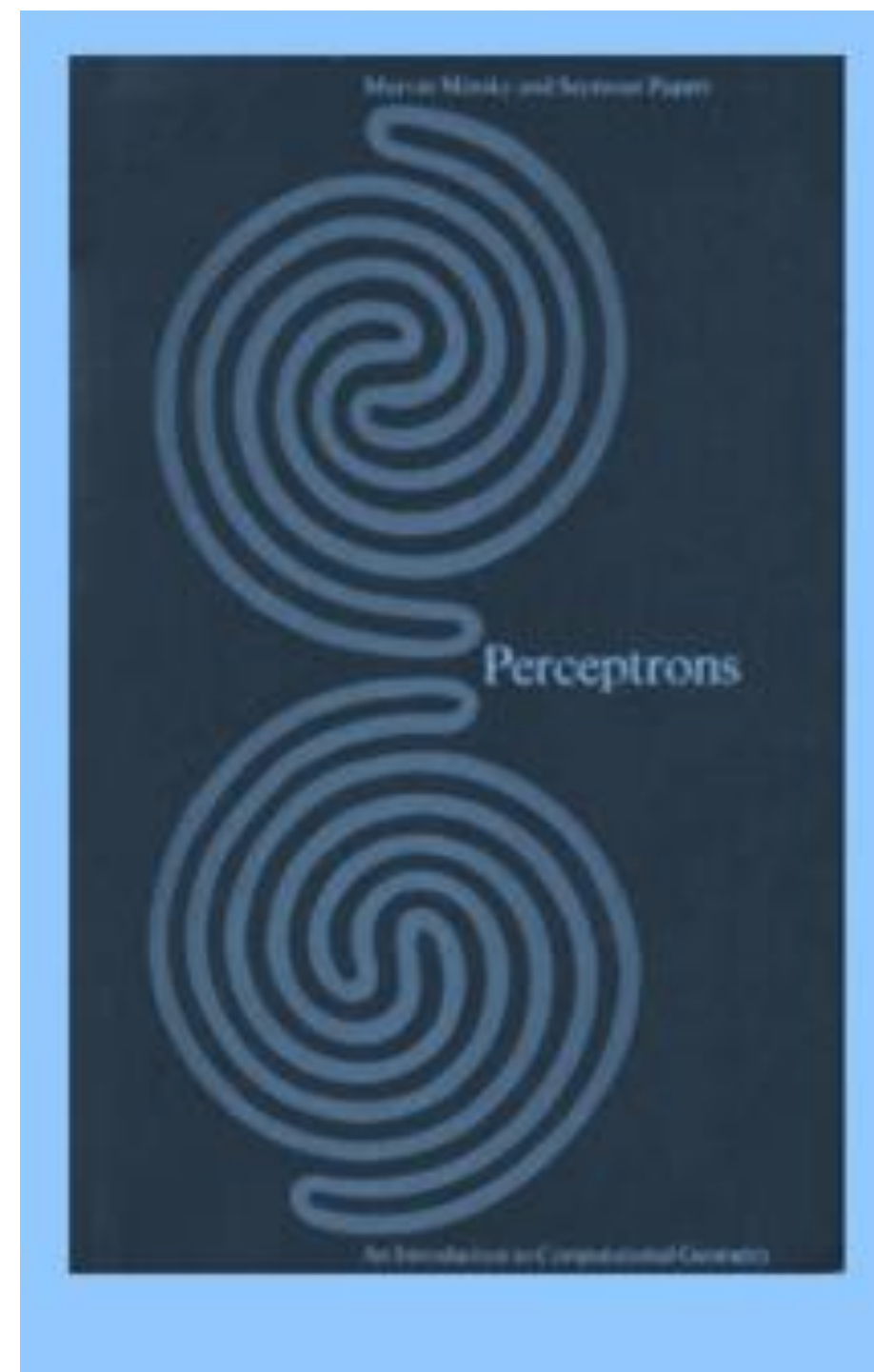
```
# cost function
cost = -tf.reduce_mean(Y*tf.log(hypothesis) + (1-Y)*tf.log(1-hypothesis))
```



# Back propagation



# Perceptrons (1969) : by Marvin Minsky, founder of the MIT AI Lab



- We need to use **MLP, multilayer perceptrons** (multilayer neural nets)
- No one on earth had found a viable way to train MLPs good enough to learn such simple functions.

