



TBSI 清华-伯克利深圳学院  
Tsinghua-Berkeley Shenzhen Institute

# Learning From Data



## Lecture 6: Backpropagation and Neural Networks

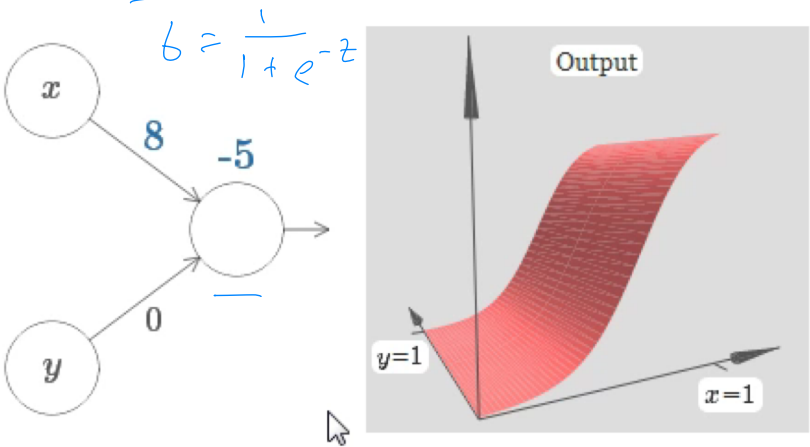
Yang Li

Slides by Lichen Wang

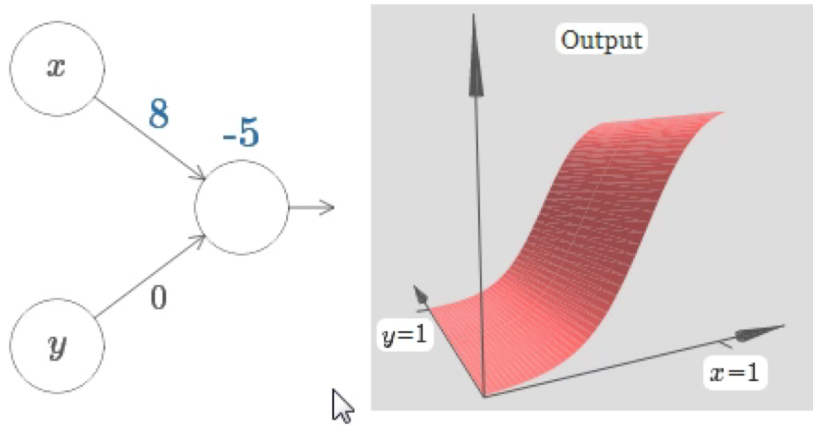
10/28/2022



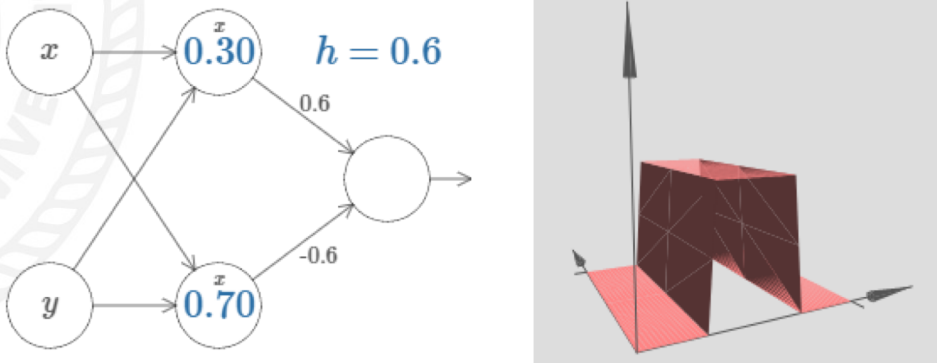
## Power of single neural



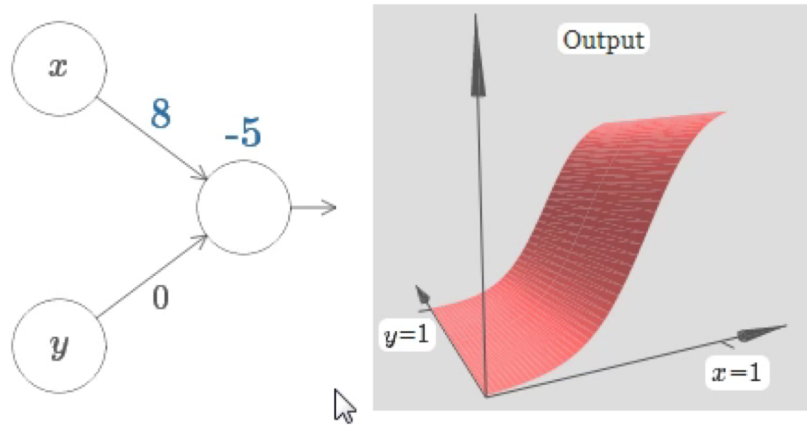
## Power of single neural



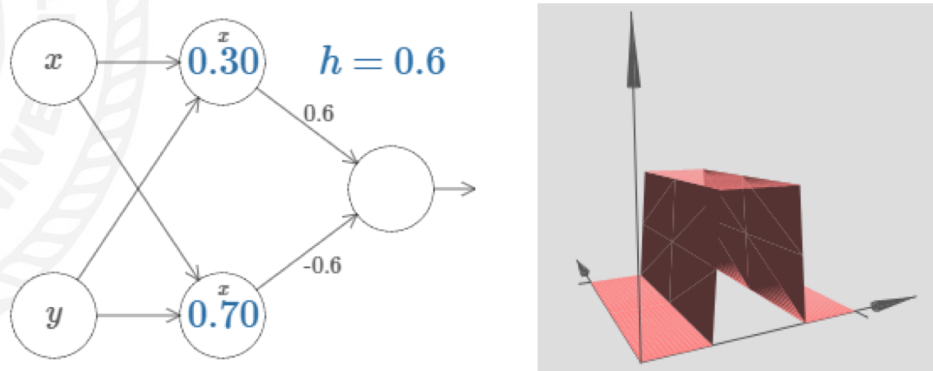
## Two hidden units



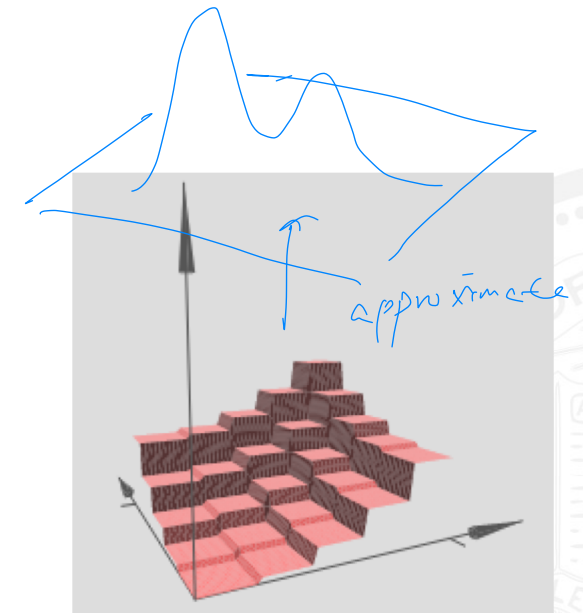
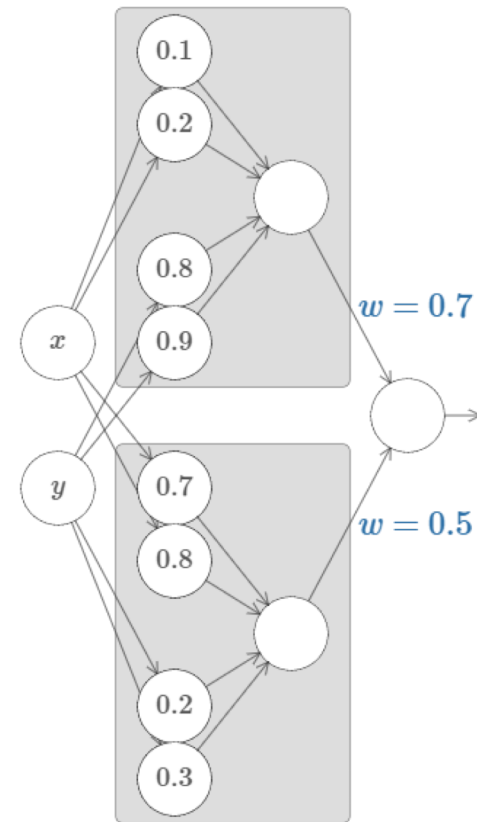
## Power of single neural



## Two hidden units



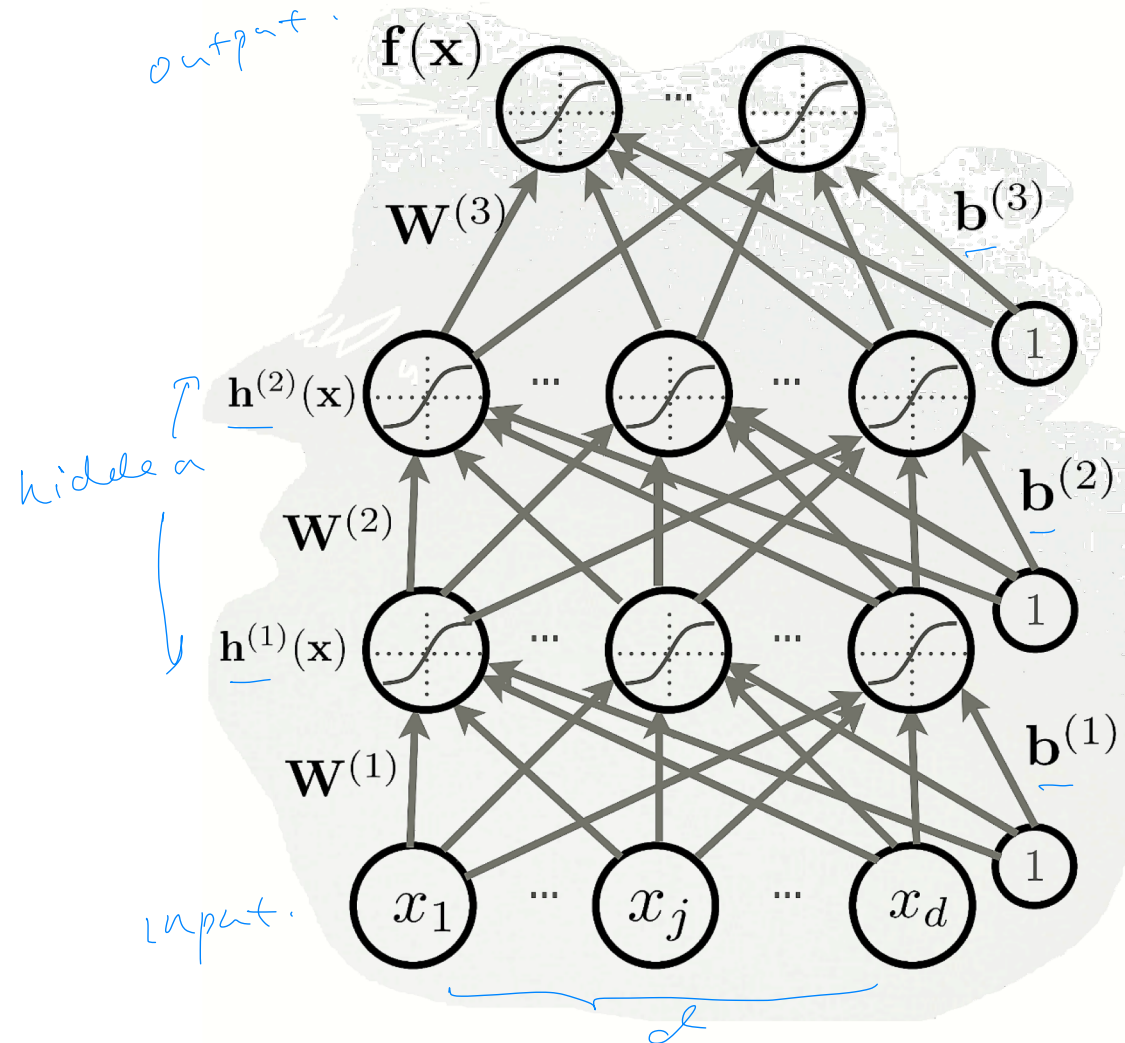
## Many hidden units





# Multilayer Neural Network

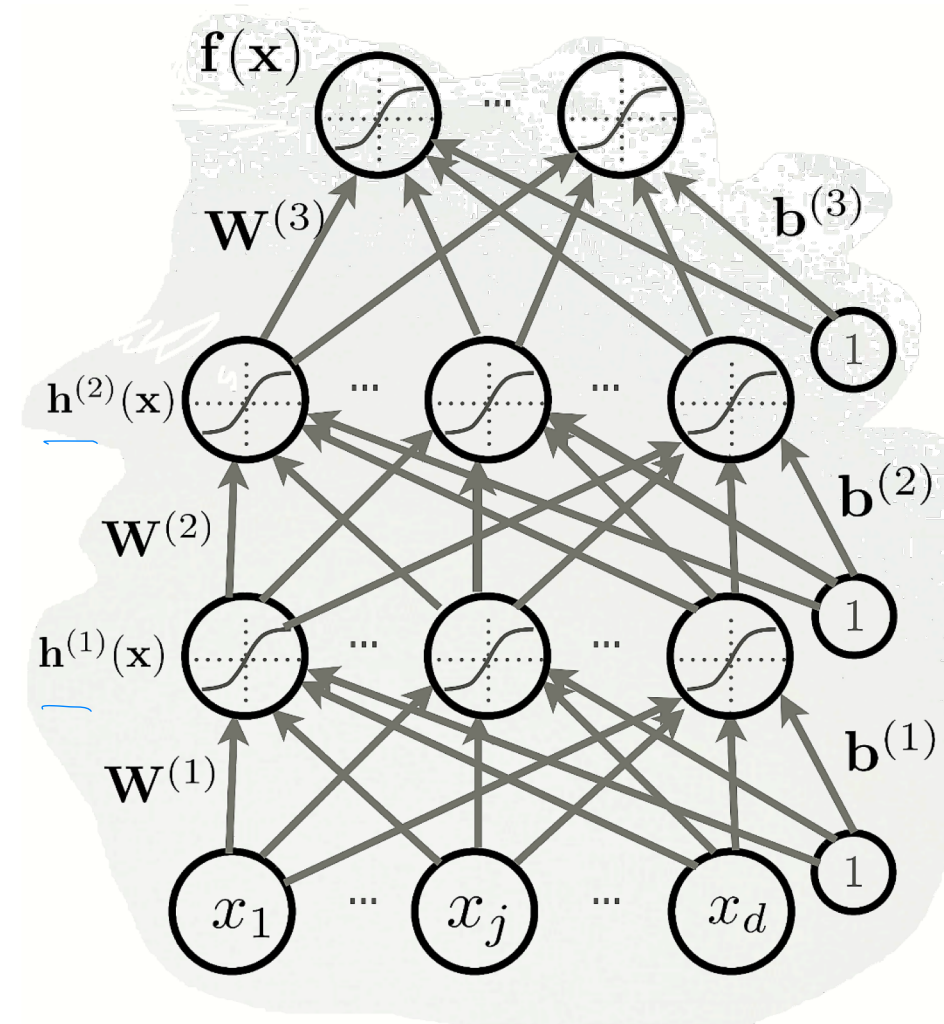
Could have L hidden layers



Could have L hidden layers

- layer input activation for  $k > 0$ ,  $\mathbf{h}^{(0)}(\mathbf{x}) = \mathbf{x}$

$$\mathbf{a}^{(k)}(\mathbf{x}) = \mathbf{b}^{(k)} + \mathbf{W}^{(k)} \mathbf{h}^{(k-1)}(\mathbf{x})$$



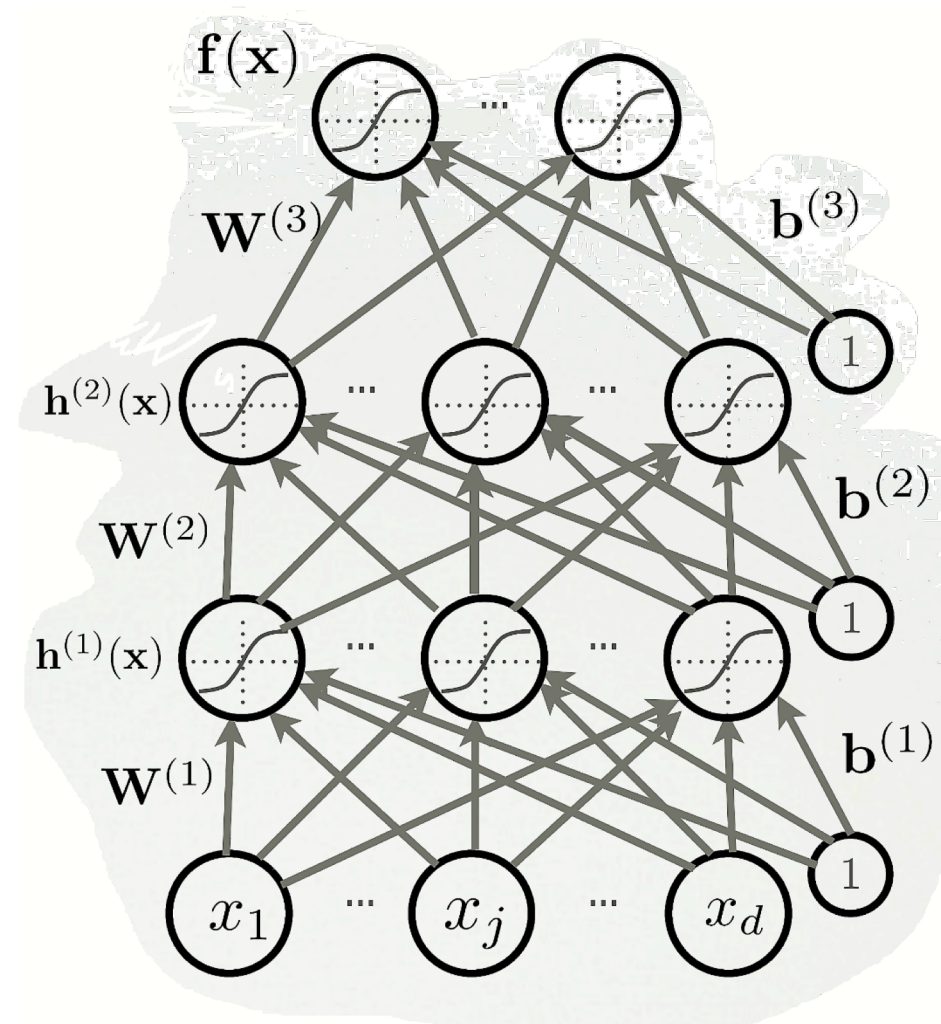
Could have  $L$  hidden layers

- layer input activation for  $k > 0$ ,  $\mathbf{h}^{(0)}(\mathbf{x}) = \mathbf{x}$

$$\mathbf{a}^{(k)}(\mathbf{x}) = \mathbf{b}^{(k)} + \mathbf{W}^{(k)} \mathbf{h}^{(k-1)}(\mathbf{x})$$

- hidden layer activation for  $1 \leq k \leq L$

$$\mathbf{h}^{(k)}(\mathbf{x}) = \mathbf{g}(\mathbf{a}^{(k)}(\mathbf{x})) \quad \text{e.g. } h^{(1)} \text{ or } h^{(2)}$$



# Multilayer Neural Network



Could have  $L$  hidden layers

- layer input activation for  $k > 0$ ,  $\mathbf{h}^{(0)}(\mathbf{x}) = \mathbf{x}$

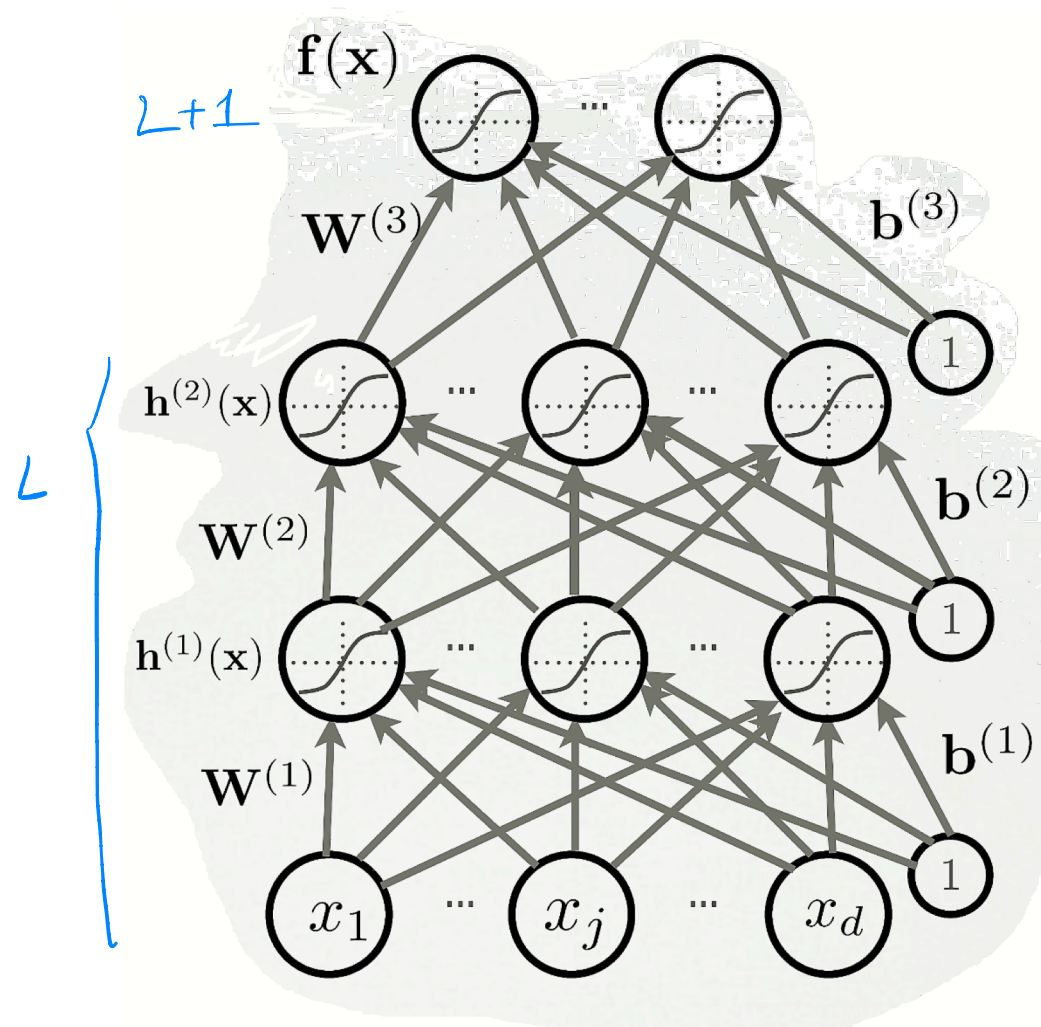
$$\mathbf{a}^{(k)}(\mathbf{x}) = \mathbf{b}^{(k)} + \mathbf{W}^{(k)} \mathbf{h}^{(k-1)}(\mathbf{x})$$

- hidden layer activation for  $1 \leq k \leq L$

$$\mathbf{h}^{(k)}(\mathbf{x}) = \mathbf{g}(\mathbf{a}^{(k)}(\mathbf{x}))$$

- output layer activation for  $k = L + 1$

$$\mathbf{h}^{(L+1)}(\mathbf{x}) = \mathbf{g}(\mathbf{a}^{(L+1)}(\mathbf{x})) = \mathbf{f}(\mathbf{x}) \text{ output.}$$





## Empirical risk

$$\operatorname{argmin}_W \frac{1}{m} \sum_{i=1}^m L(f(\mathbf{x}^{(i)}; W), y^{(i)}) + \lambda \Omega(W)$$

■  $L(\underbrace{f(\mathbf{x}^{(i)}; W)}_{\hat{y}^{(i)}}, y^{(i)})$  is the loss function for sample  $(\mathbf{x}^{(i)}, y^{(i)})$

■  $\lambda \Omega(W)$  is the regularizer

$$p(y=i|x) = \frac{e^{\theta_i^T x}}{\sum_j e^{\theta_j^T x}}$$

e. g. When L is the softmax loss

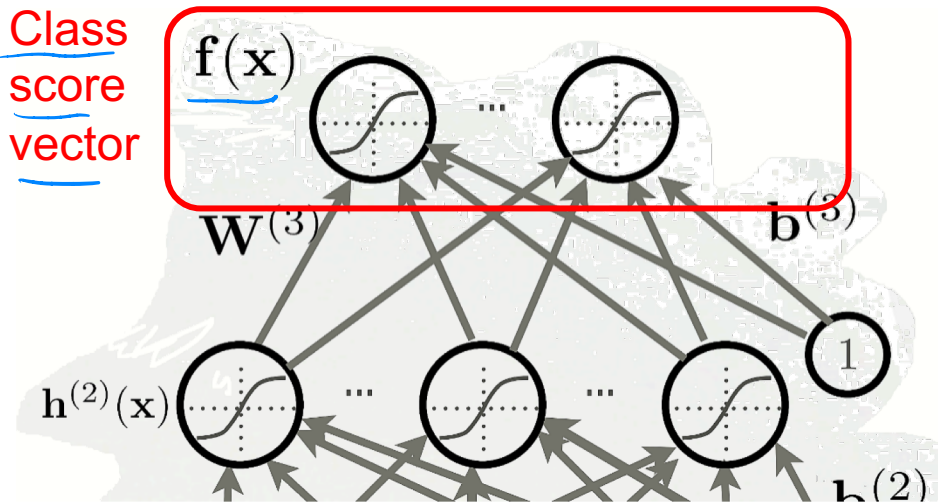
$$L(f(\mathbf{x}^{(i)}; W), y^{(i)}) = -\log \left( \frac{e^{f_{y^{(i)}}}}{\sum_{j=1}^{|Y|} e^{f_j}} \right)$$

$f_j$  is the  $j$ th element of class score vector  $f(\mathbf{x}^{(i)}; W)$

Softmax example:

Unnormalized class probability of  $|Y|$  classes

Class score vector





- Find the optimal parameter

$$\operatorname{argmin}_W \frac{1}{m} \sum_{i=1}^m L(f(\mathbf{x}^{(i)}; W), y^{(i)}) + \lambda \Omega(W)$$



## ■ Find the optimal parameter

$$\operatorname{argmin}_W \frac{1}{m} \sum_{i=1}^m L(f(\mathbf{x}^{(i)}; W), y^{(i)}) + \lambda \Omega(W)$$

To apply this algorithm, we need:

1. A procedure to compute the parameter gradient
2. The regularizer (and its gradient)
3. Updating rule
4. Initialization method

## ■ Stochastic Gradient Descent (SGD)

### Algorithm

#### 1. Initialize W

repeat: for each training example  $(\mathbf{x}^{(t)}, y^{(t)})$

2a.  $\Delta = -\nabla_{\mathbf{w}} L(f(\mathbf{x}^{(t)}; \mathbf{W}), y^{(t)}) - \lambda \nabla_{\mathbf{w}} \Omega(\mathbf{W})$

2b.  $\mathbf{W} \leftarrow \mathbf{W} + \alpha \Delta$

Training epoch

=

Iterating over all examples

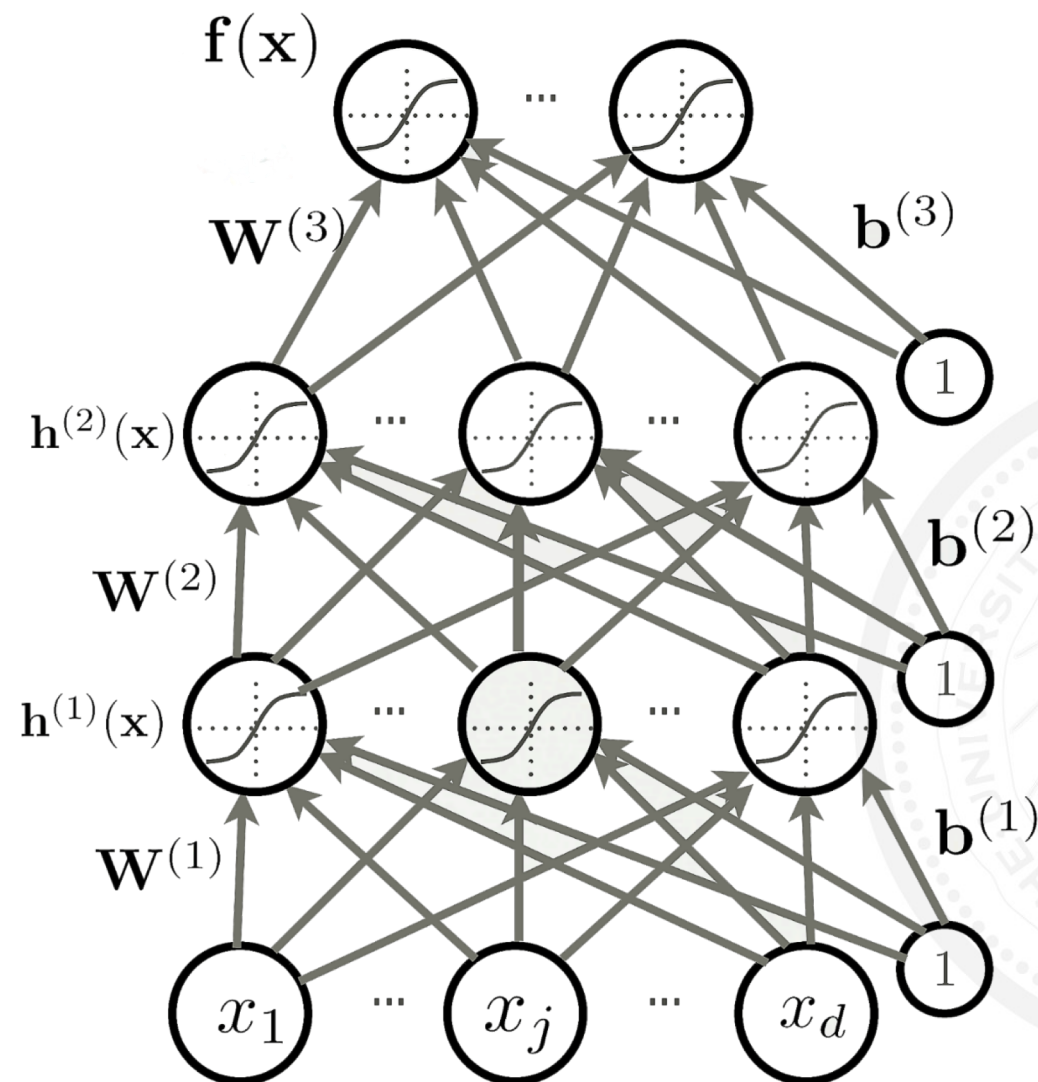
Animation:

<http://vision.stanford.edu/teaching/cs231n-demos/linear-classify/>

# Follow the slope



How many parameter do we have?

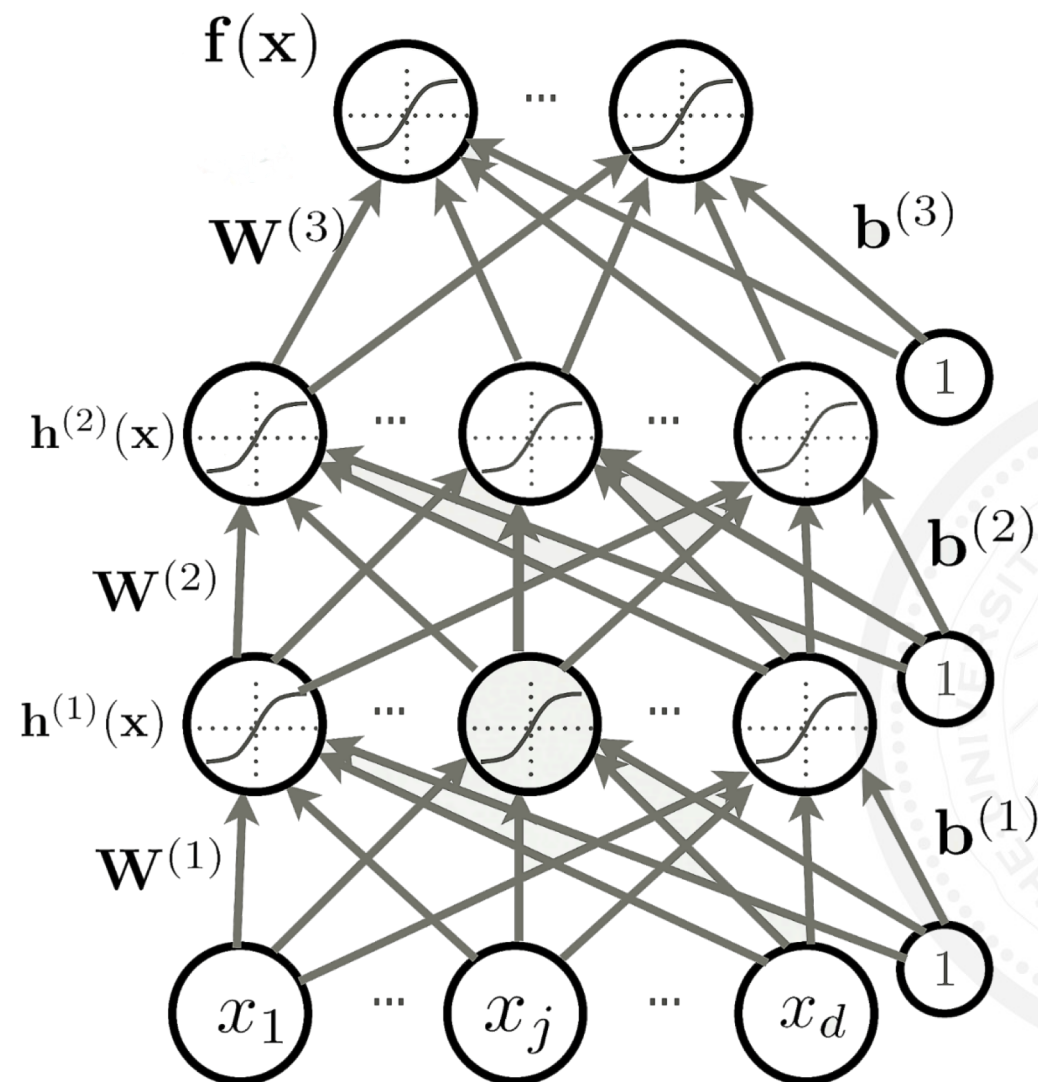




# Follow the slope

How many parameter do we have?

VGGNet [Simonyan and Zisserman, 2014] used 138M parameters



# Follow the slope

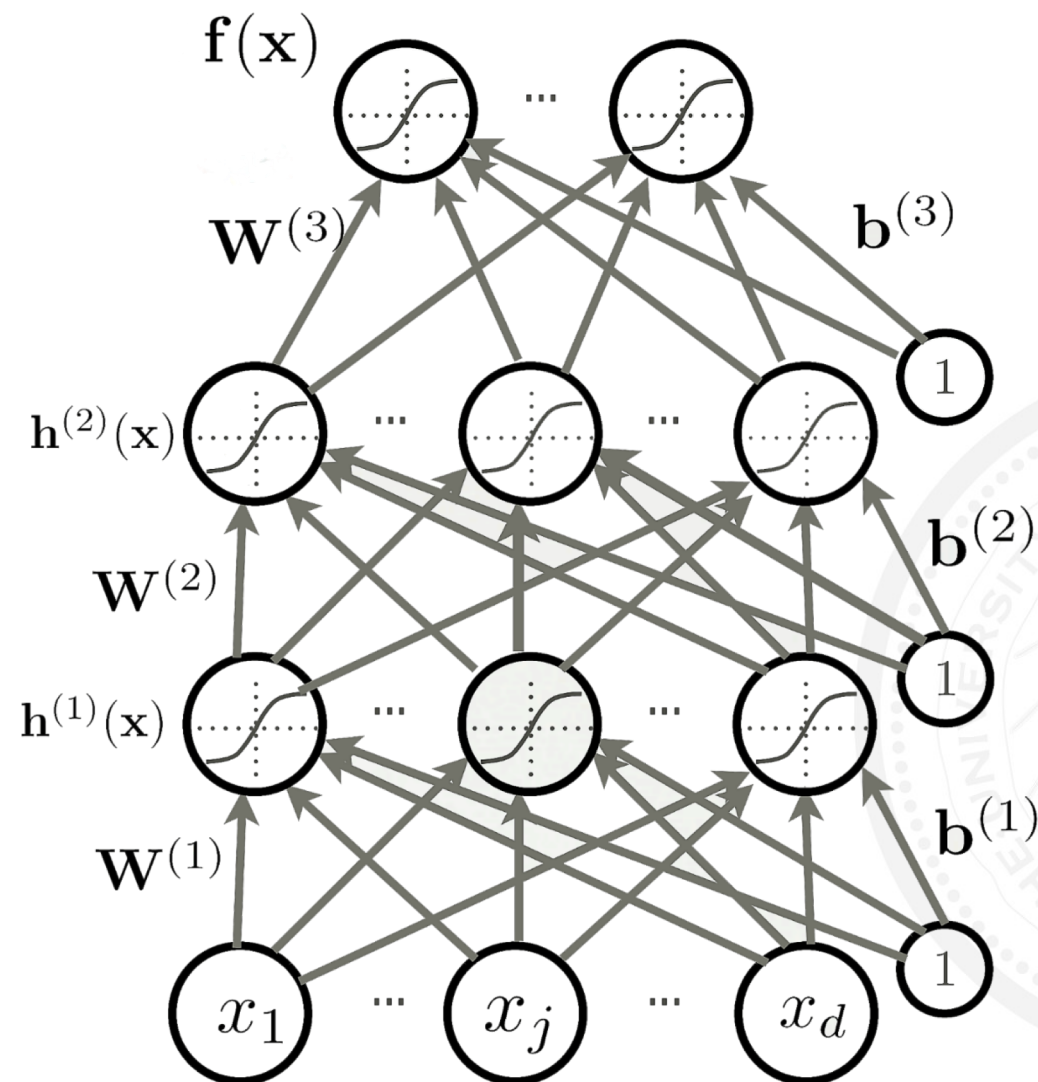


How many parameter do we have?

VGGNet [Simonyan and Zisserman, 2014] used 138M parameters

In 1-dimension, the derivative of a function:

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



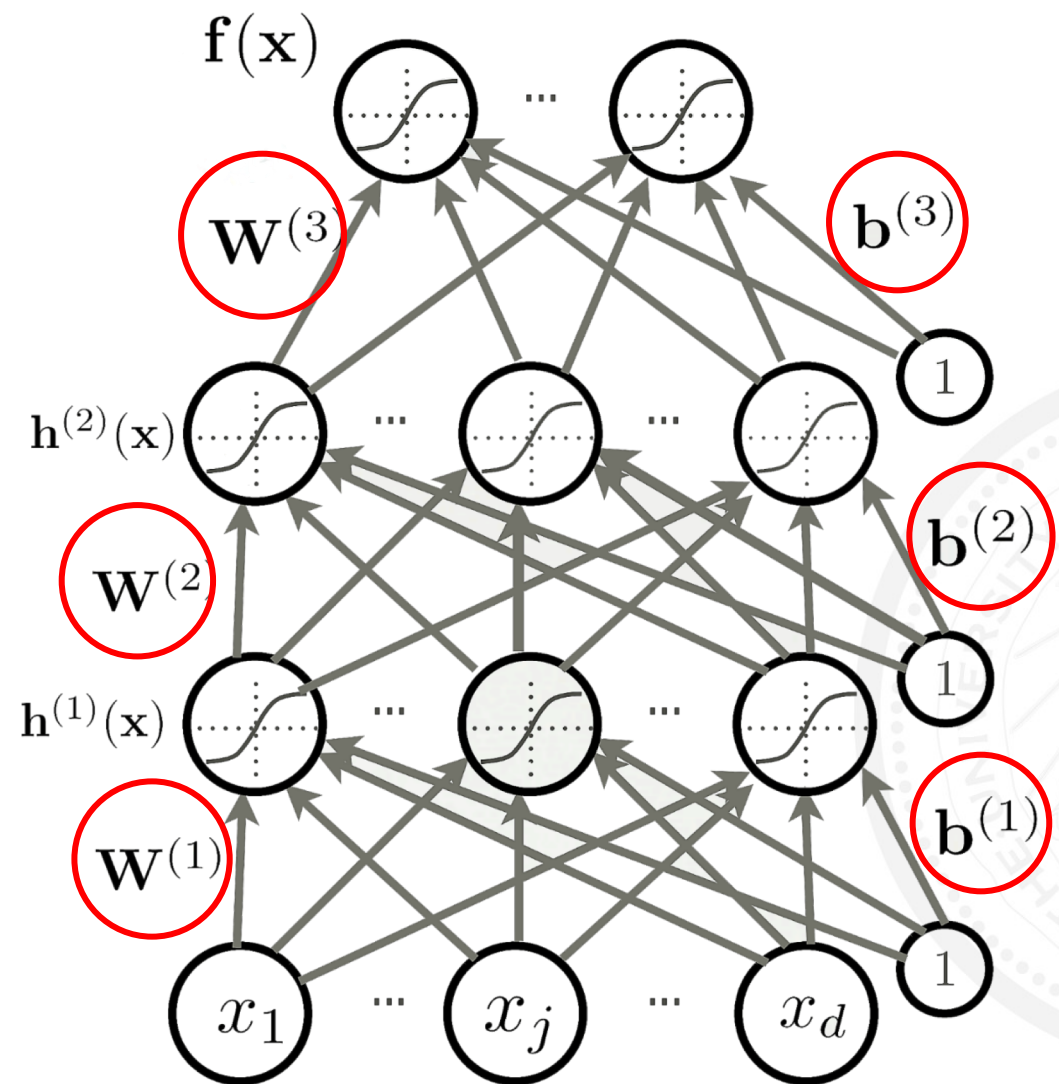
# Follow the slope

How many parameter do we have?

VGGNet [Simonyan and Zisserman, 2014] used 138M parameters

In 1-dimension, the derivative of a function:

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



# Numerical Gradient



TBSI 清华-伯克利深圳学院  
Tsinghua-Berkeley Shenzhen Institute

Current  $W$ :

Gradient  $dW$ :



Current W:

[0.25,  
-1.56,  
0.55,  
3.8,  
0.98,  
0.77,  
-0.11,  
-2.9,...]

**Loss 1.25742**

$$\partial w = \lim_{h \rightarrow 0} \frac{f(w+h) - f(w)}{h}$$

Gradient dW:

[ ?,  
?,  
?,  
?,  
?,  
?,  
?,  
?,  
?,...]



Current  $W$ :

[0.25,  
-1.56,  
0.55,  
3.8,  
0.98,  
0.77,  
-0.11,  
-2.9,...]

**Loss 1.25742**

$f(w)$

$W$  +  $h$  (third dim):

[0.25 + 0.0001,  
-1.56,  
0.55,  
3.8,  
0.98,  
0.77,  
-0.11,  
-2.9,...]

**Loss 1.25763**

$f(w+h)$

Gradient  $dW$ :

[ ? ,  
 ? ,  
 ? ,  
 ? ,  
 ? ,  
 ? ,  
 ? ,  
 ? , ... ]

Current W:

[0.25,  
-1.56,  
0.55,  
3.8,  
0.98,  
0.77,  
-0.11,  
-2.9,...]

**Loss 1.25742**

W + h (third dim):

$W_{01}$  [0.25 + 0.0001,  
-1.56,  
0.55,  
3.8,  
0.98,  
0.77,  
-0.11,  
-2.9,...]

**Loss 1.25763**

Gradient dW:

[2.1,  
?,  
?,  
?,  
?,  
?,  
?,  
?,  
?,...]

$$= \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} = \frac{(1.25763 - 1.25742)}{0.0001}$$



# Follow the slope

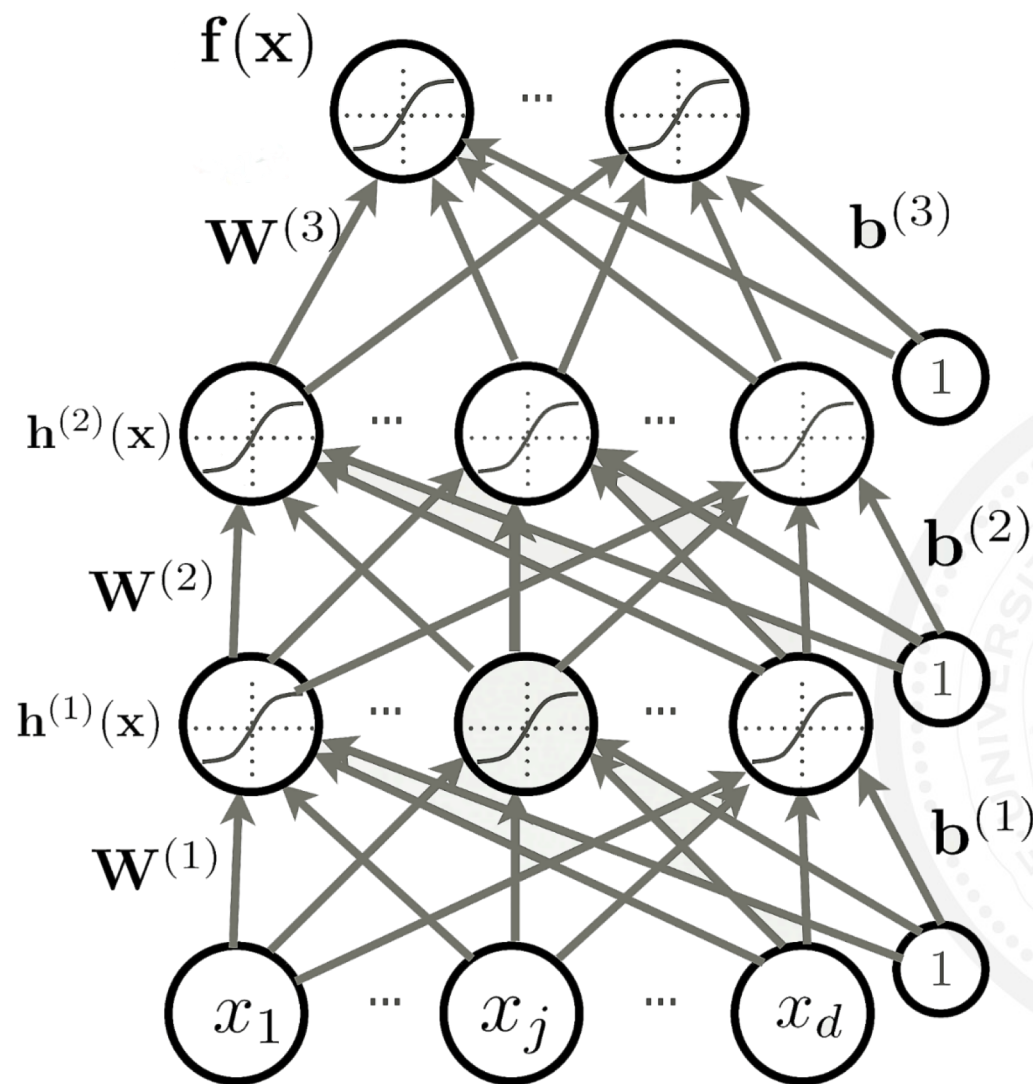


How many parameter do we have?

VGGNet [Simonyan and Zisserman, 2014] used 138M parameters

In 1-dimension, the derivative of a function:

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





# Follow the slope



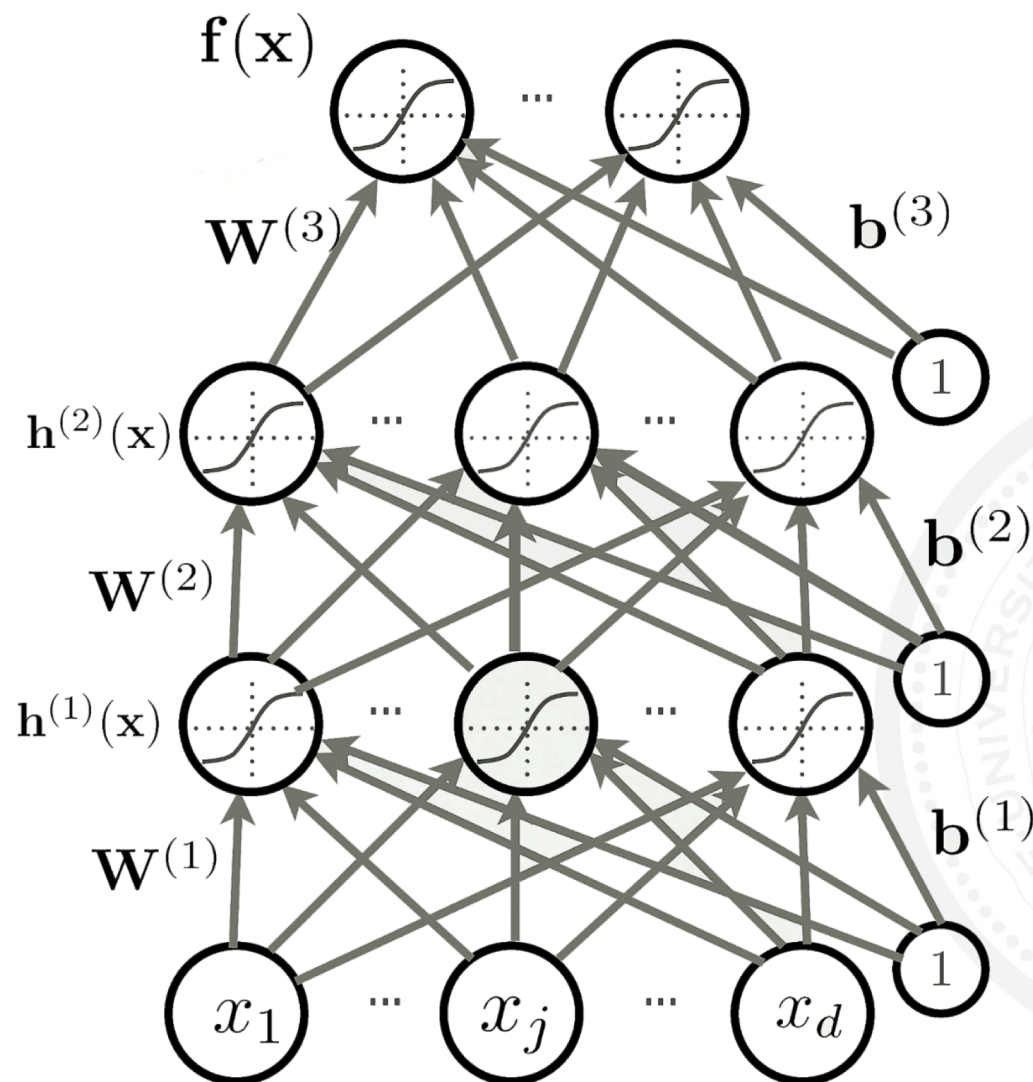
How many parameter do we have?

VGGNet [Simonyan and Zisserman, 2014] used 138M parameters

In 1-dimension, the derivative of a function:

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

Numerical gradient: approximate, slow, easy to write



# Follow the slope

How many parameter do we have?

VGGNet [Simonyan and Zisserman, 2014] used 138M parameters

In 1-dimension, the derivative of a function:

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

Numerical gradient: approximate, slow, easy to write

Calculus!

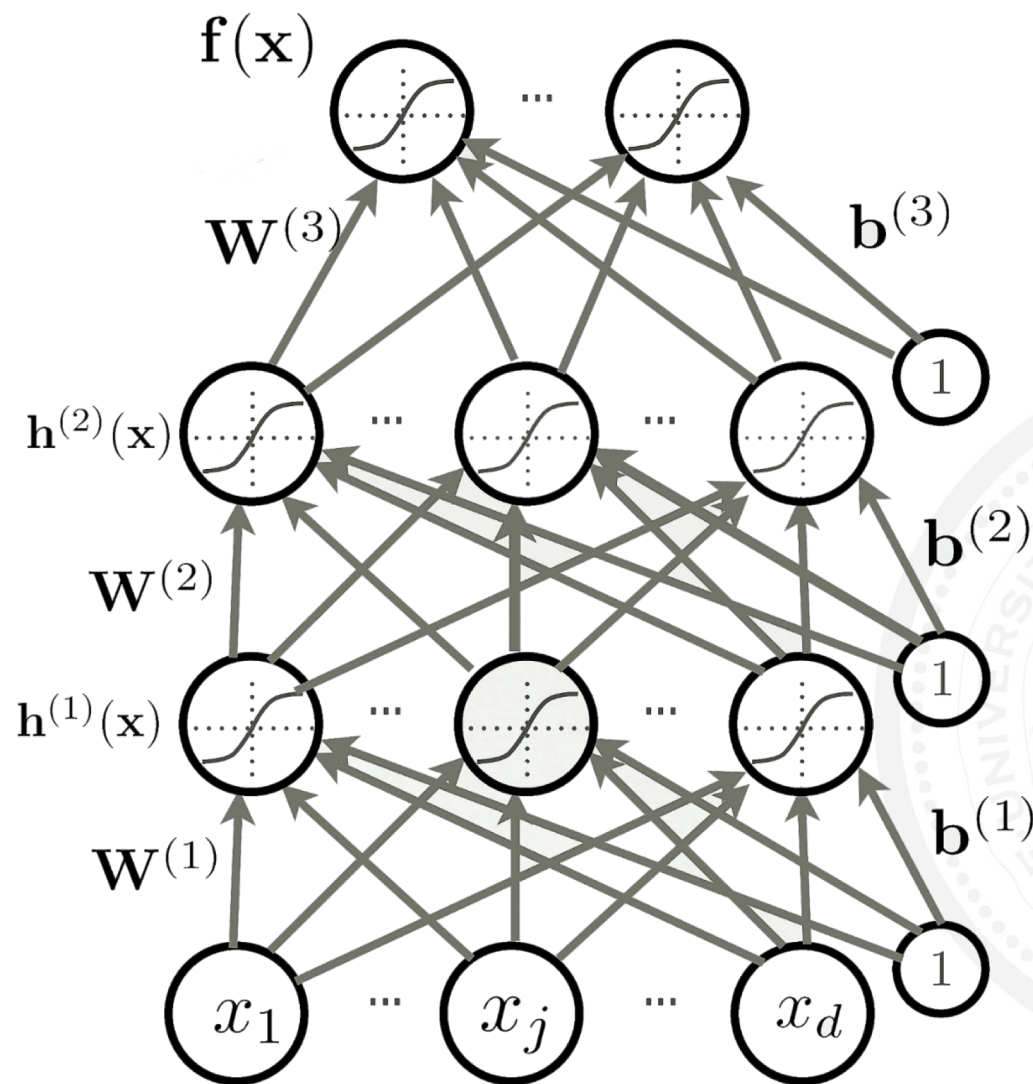
$$\operatorname{argmin}_W \frac{1}{m} \sum_{i=1}^m L(f(\mathbf{x}^{(i)}; W), y^{(i)}) + \lambda \Omega(W)$$

$$\mathbf{h}^{(L+1)}(\mathbf{x}) = \mathbf{g}(\mathbf{a}^{(L+1)}(\mathbf{x})) = \mathbf{f}(\mathbf{x})$$

$$\mathbf{h}^{(k)}(\mathbf{x}) = \mathbf{g}(\mathbf{a}^{(k)}(\mathbf{x}))$$

$$\mathbf{a}^{(k)}(\mathbf{x}) = \mathbf{b}^{(k)} + \mathbf{W}^{(k)} \mathbf{h}^{(k-1)}(\mathbf{x})$$

Analytic gradient: exact, fast, error-prone



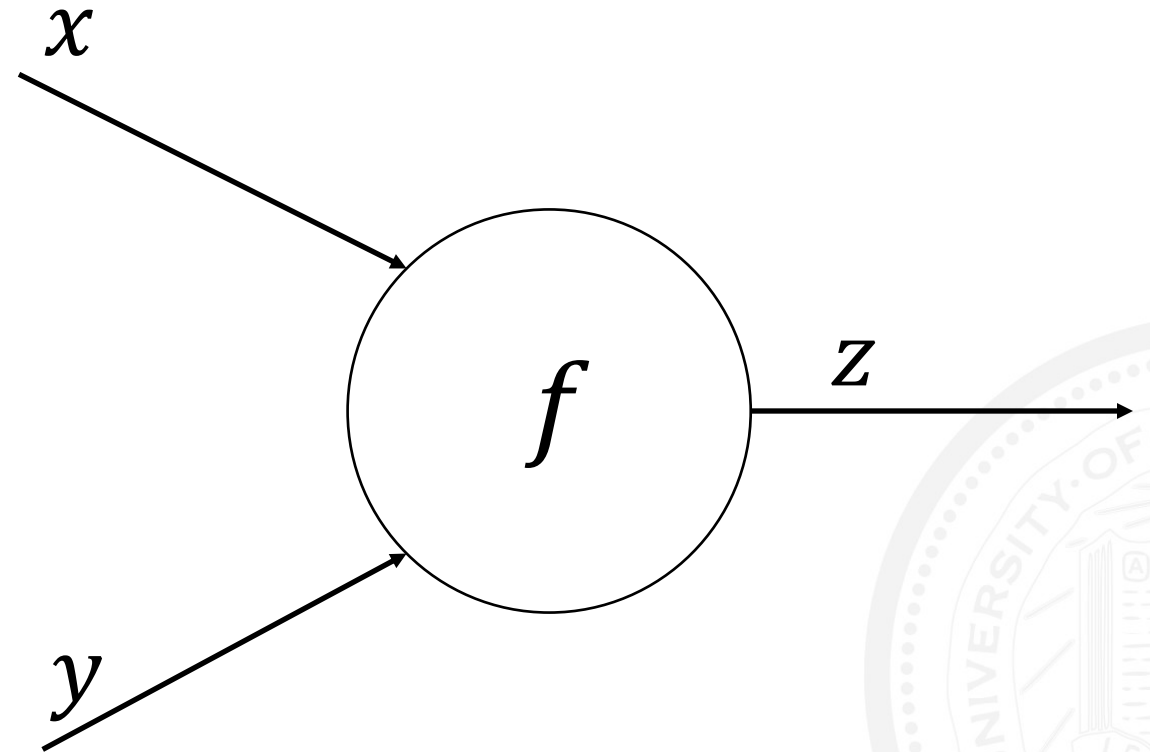
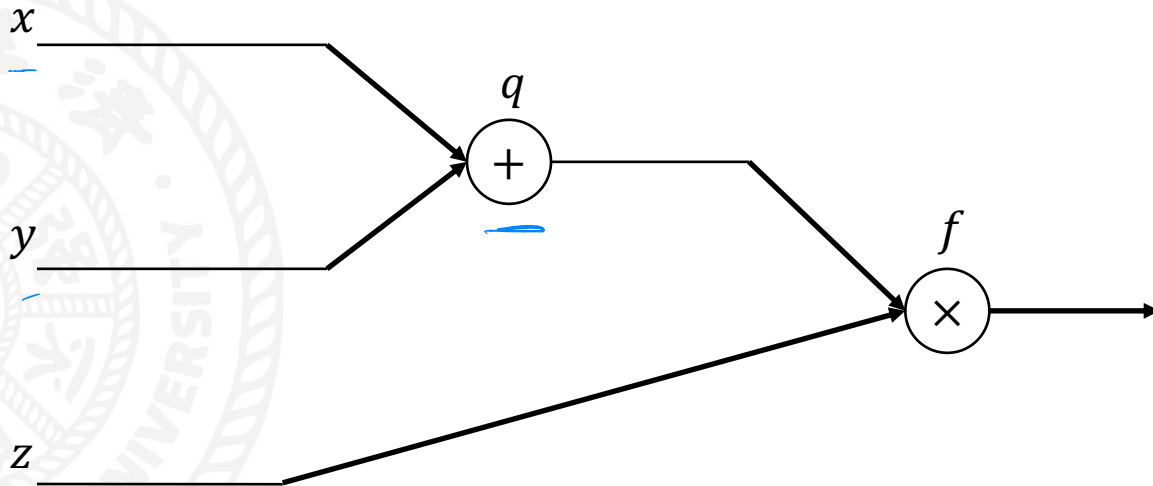
# Backpropagation



computation graph

$$f(x, y, z) = \underline{(x + y)}z$$

We want  $\frac{df}{dx}$ ,  $\frac{df}{dy}$ ,  $\frac{df}{dz}$



# Backpropagation



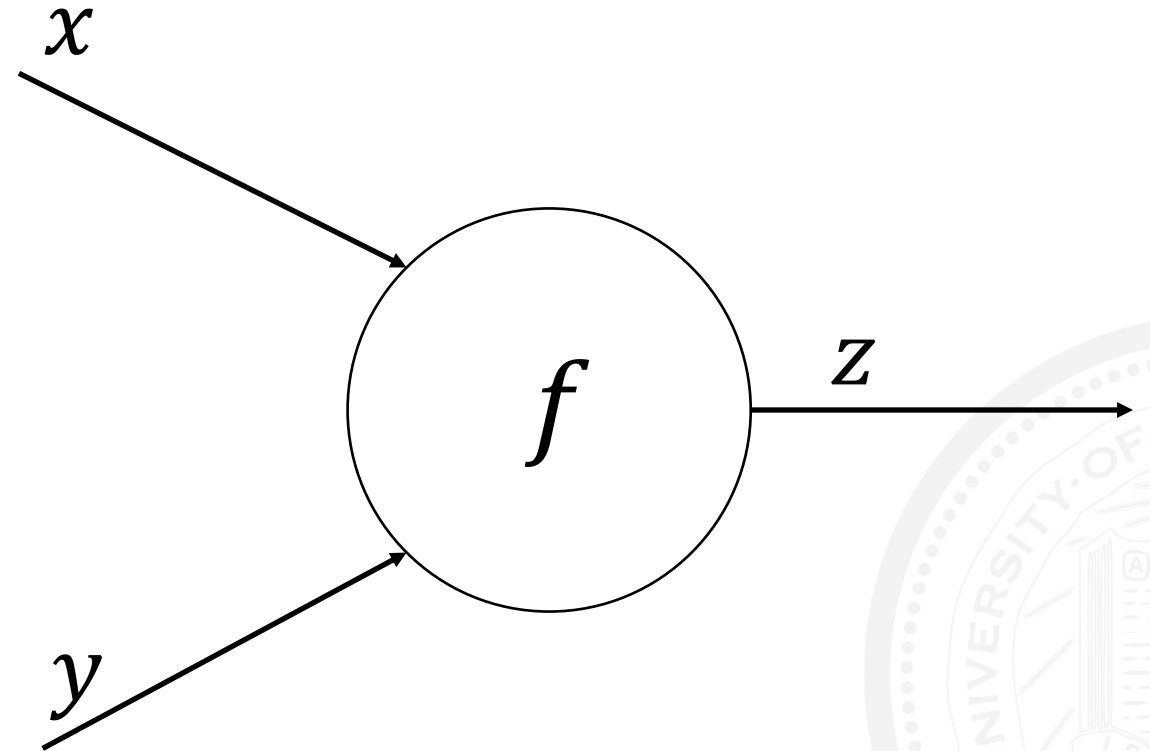
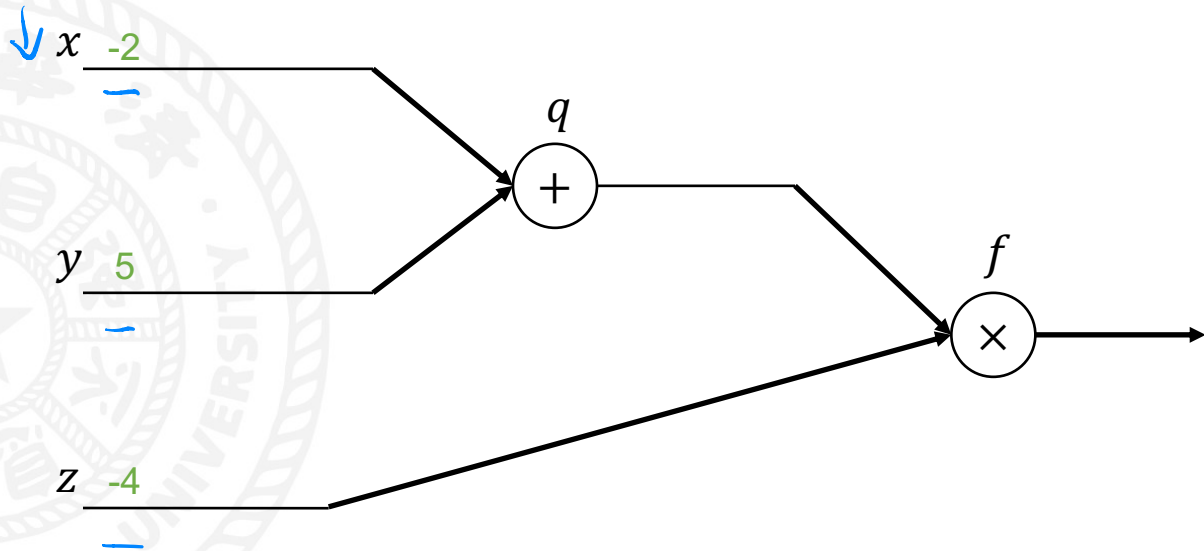
TBSI

清华-伯克利深圳学院  
Tsinghua-Berkeley Shenzhen Institute

$$f(x, y, z) = (x + y)z$$

We want  $\frac{df}{dx}$ ,  $\frac{df}{dy}$ ,  $\frac{df}{dz}$

Initialize



# Backpropagation

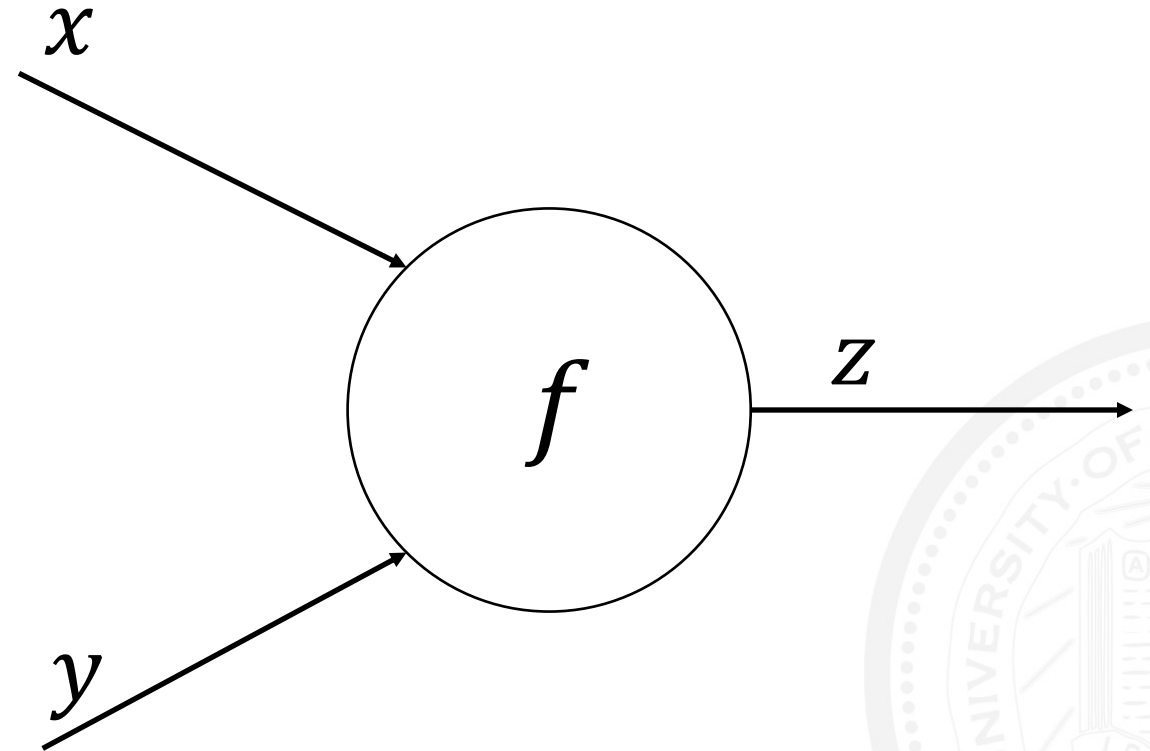
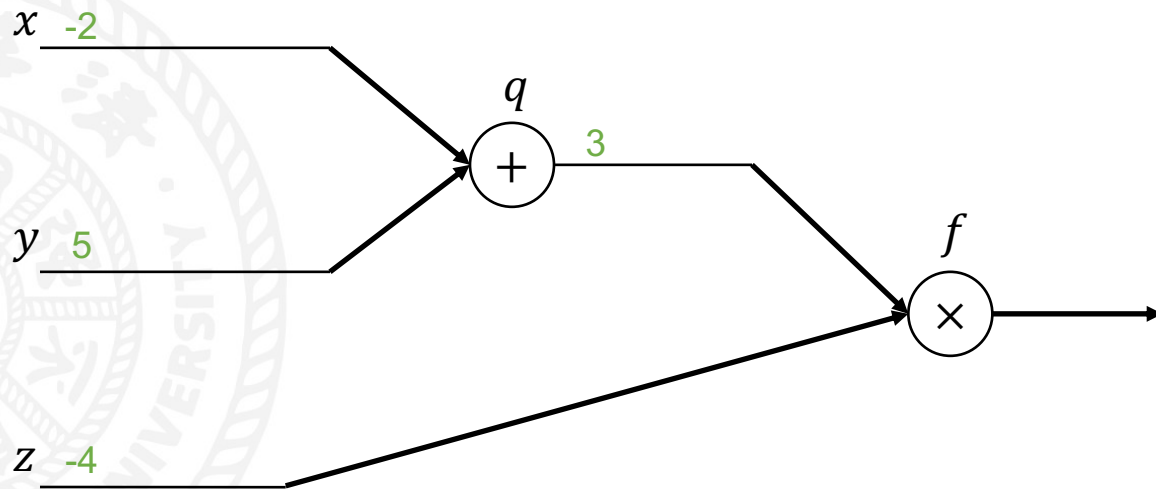


TBSI

清华-伯克利深圳学院  
Tsinghua-Berkeley Shenzhen Institute

$$f(x, y, z) = (x + y)z$$

We want  $\frac{df}{dx}$ ,  $\frac{df}{dy}$ ,  $\frac{df}{dz}$



# Backpropagation

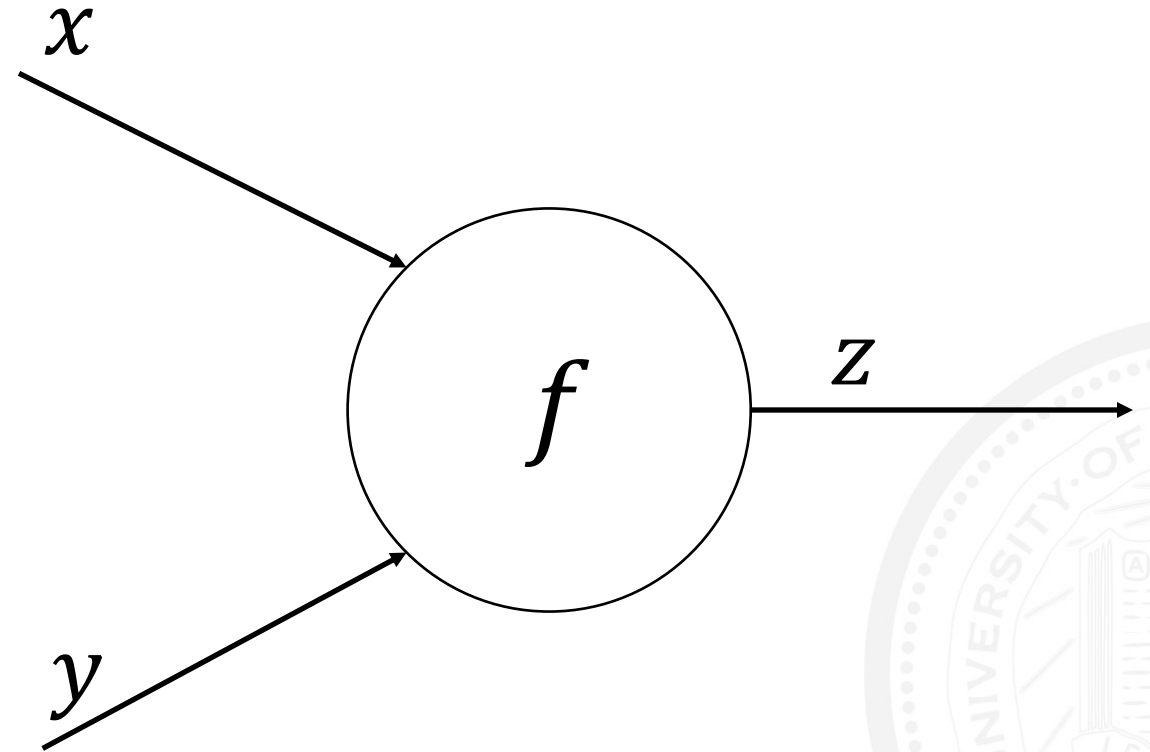
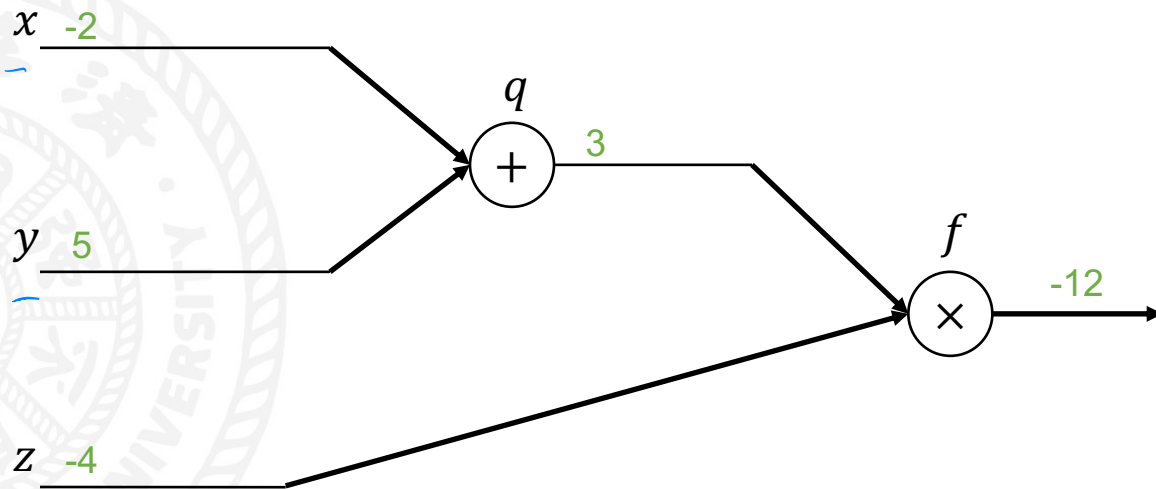


TBSI

清华-伯克利深圳学院  
Tsinghua-Berkeley Shenzhen Institute

$$f(x, y, z) = (x + y)z$$

We want  $\frac{df}{dx}$ ,  $\frac{df}{dy}$ ,  $\frac{df}{dz}$



# Backpropagation

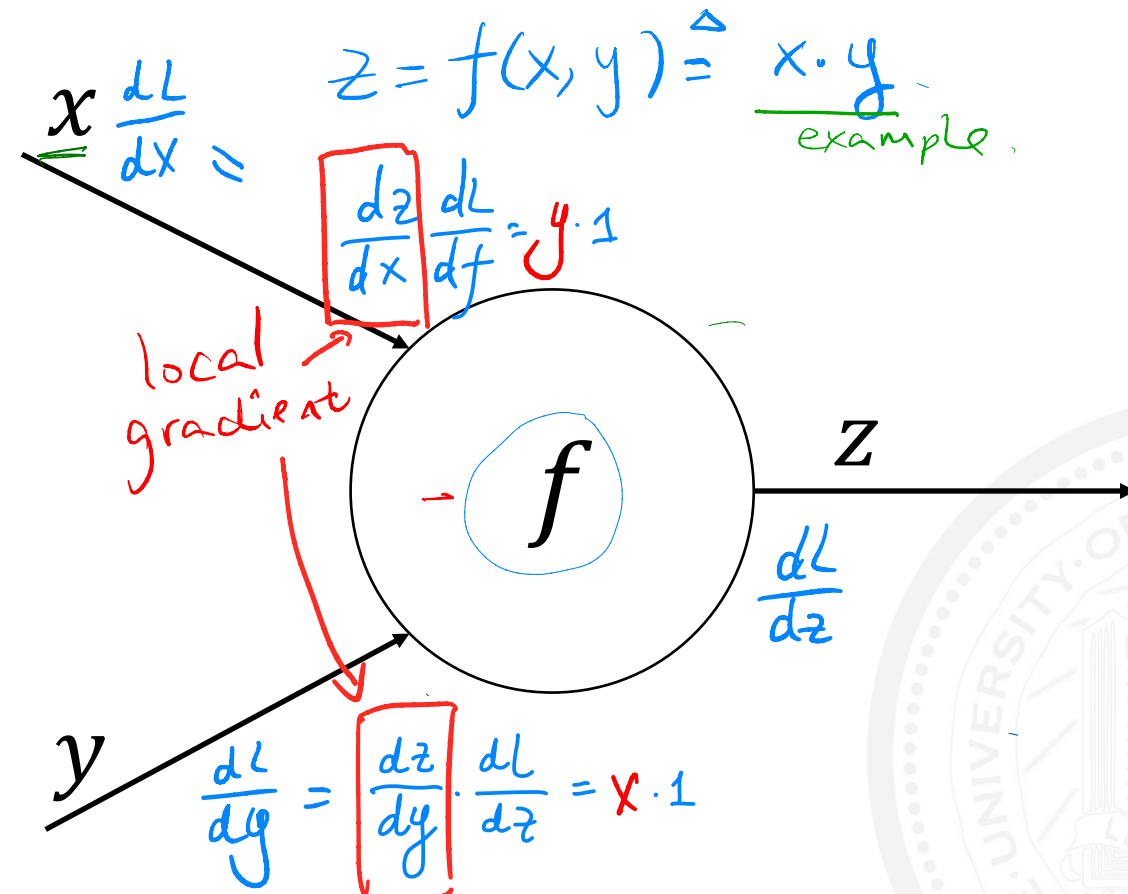
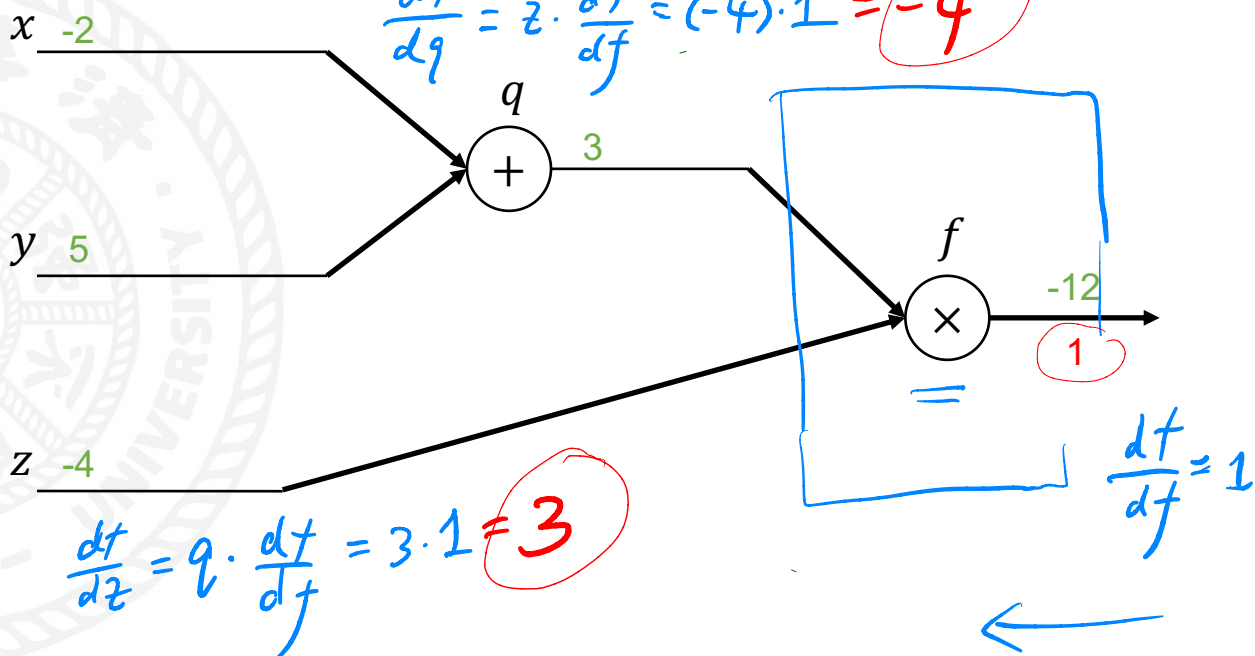


$$f(x, y, z) = (x + y)z$$

We want  $\frac{df}{dx}, \frac{df}{dy}, \frac{df}{dz}$

$$\frac{df}{dz} = z \cdot \frac{df}{df} = (-4) \cdot 1 = -4$$

$$\frac{df}{dx} = q \cdot \frac{df}{df} = 3 \cdot 1 = 3$$



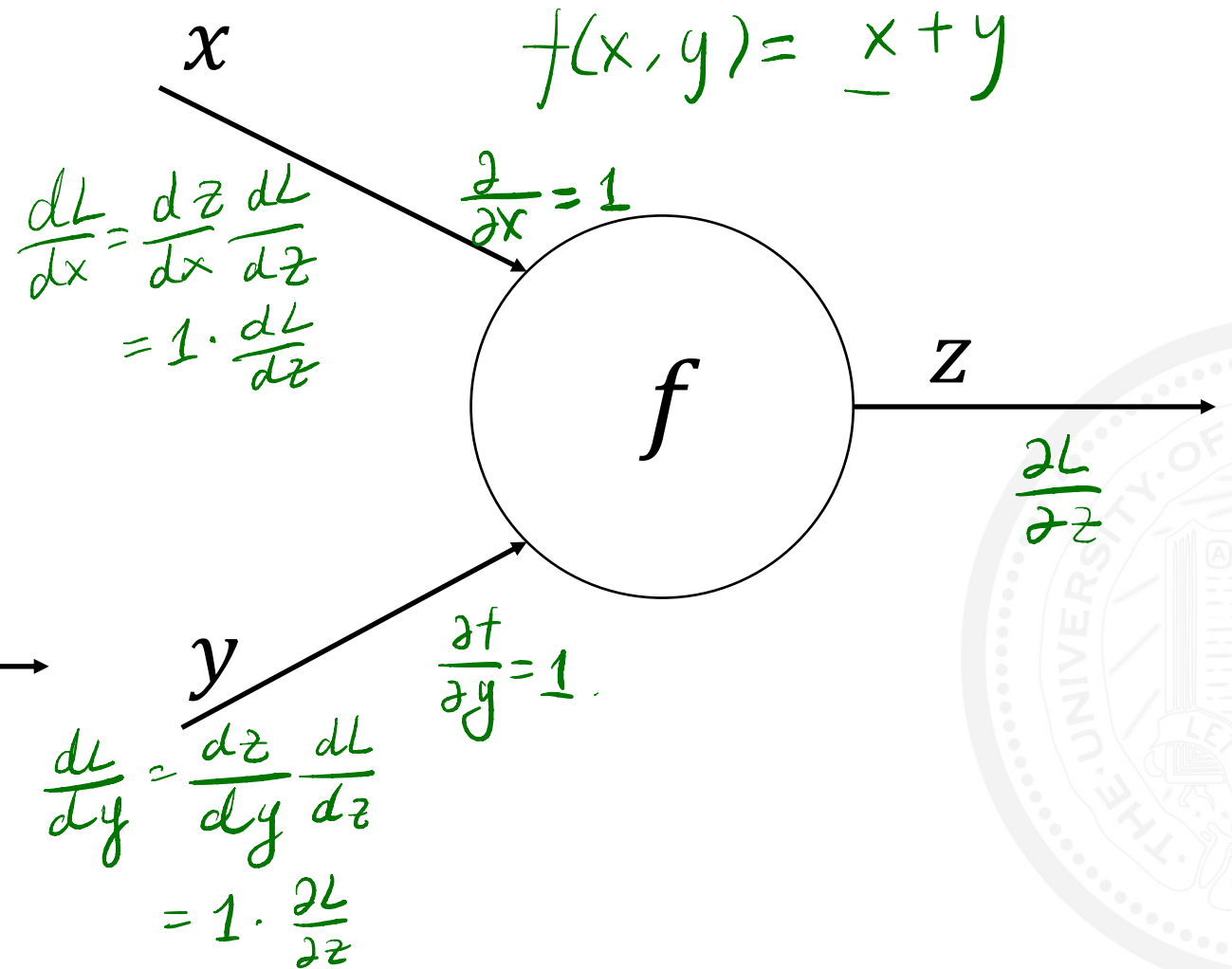
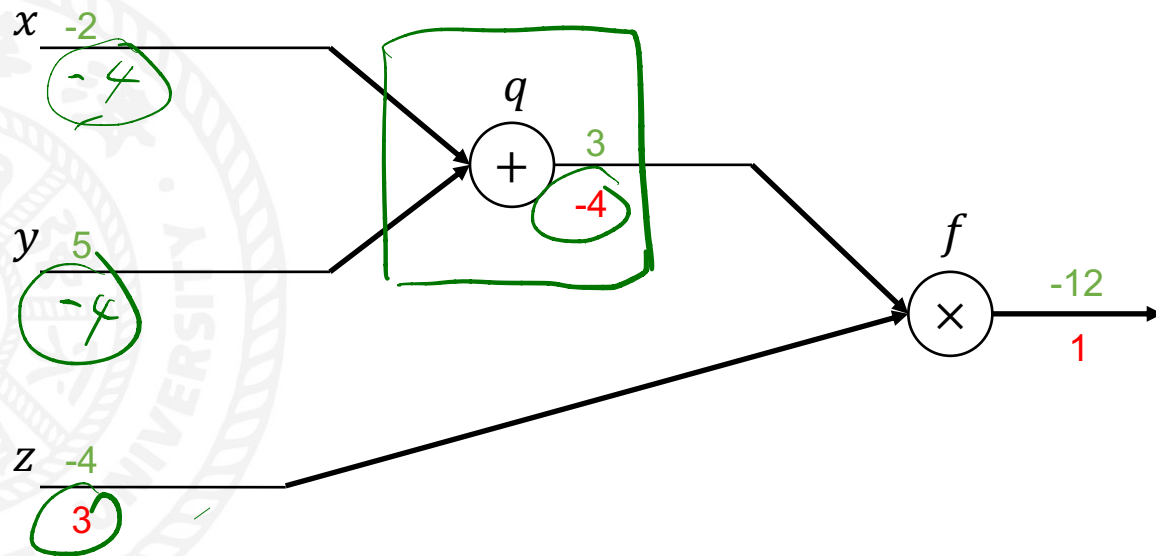


# Backpropagation



$$f(x, y, z) = (x + y)z$$

We want  $\frac{df}{dx}$ ,  $\frac{df}{dy}$ ,  $\frac{df}{dz}$



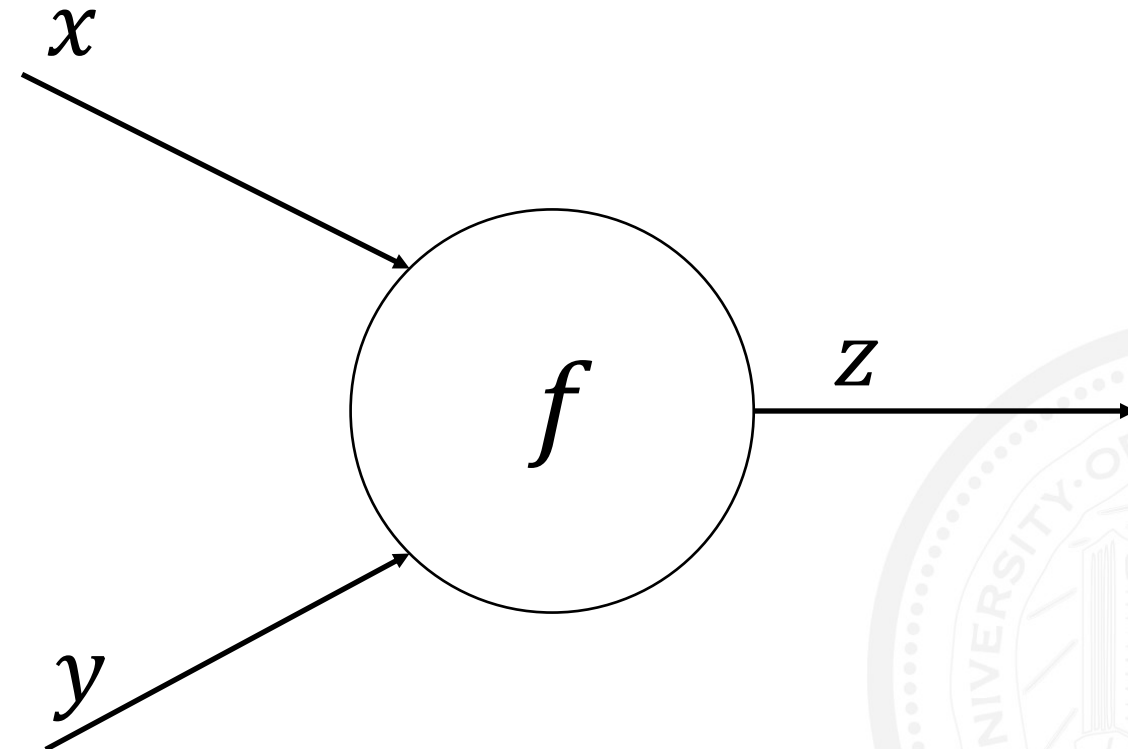
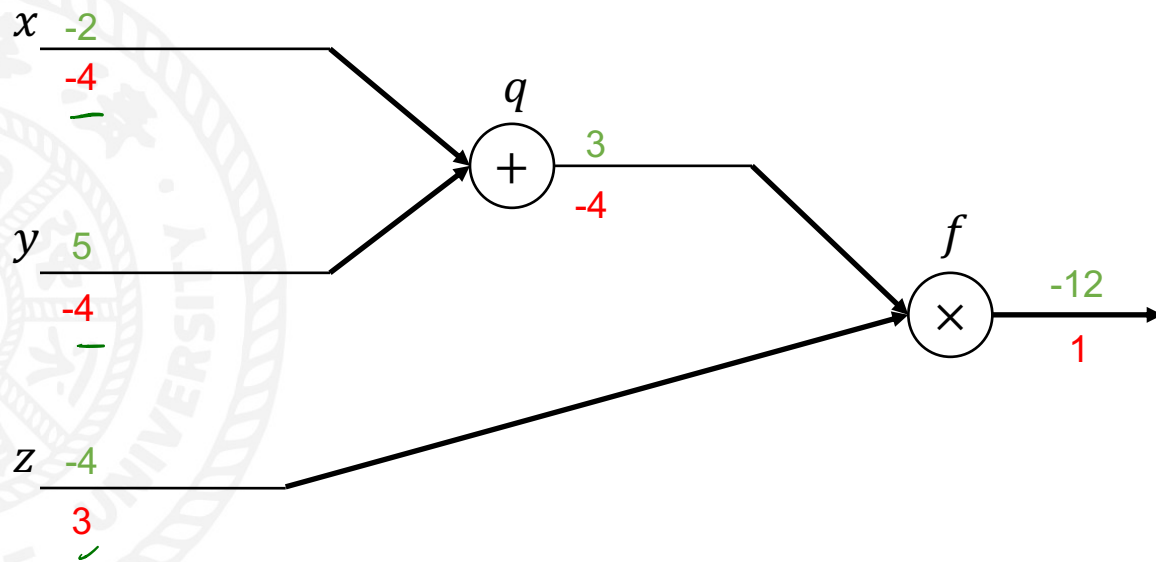


# Backpropagation



$$f(x, y, z) = (x + y)z$$

We want  $\frac{df}{dx}$ ,  $\frac{df}{dy}$ ,  $\frac{df}{dz}$

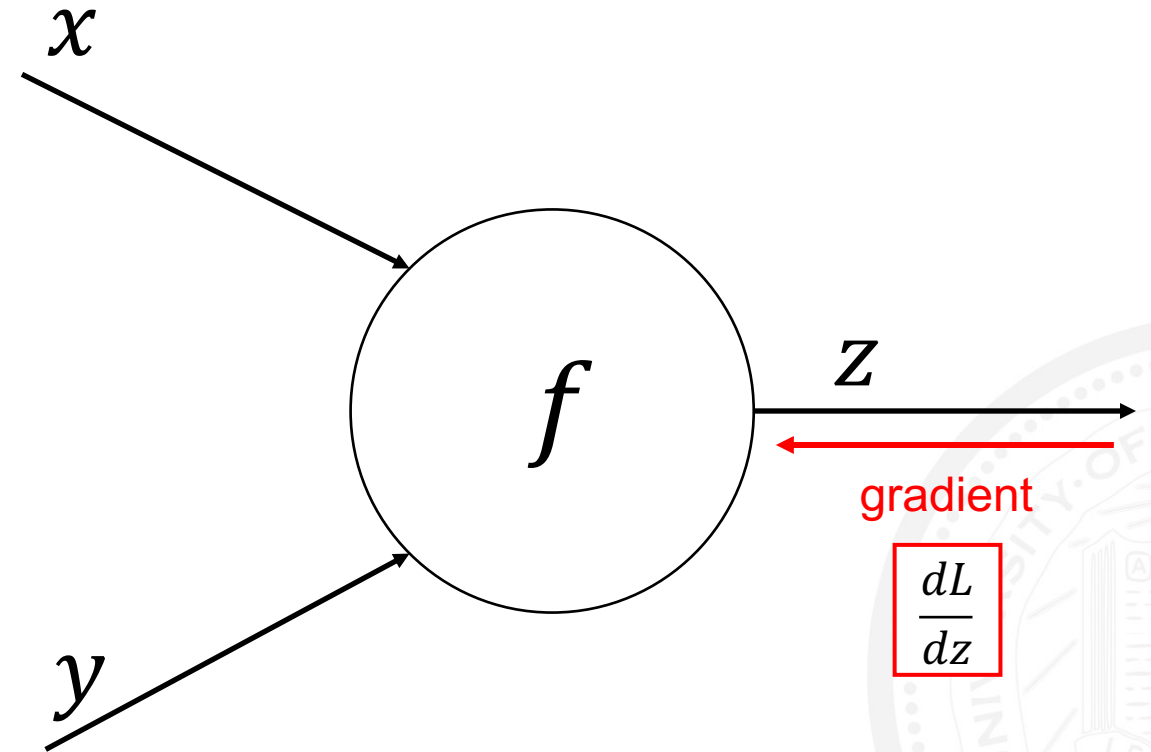
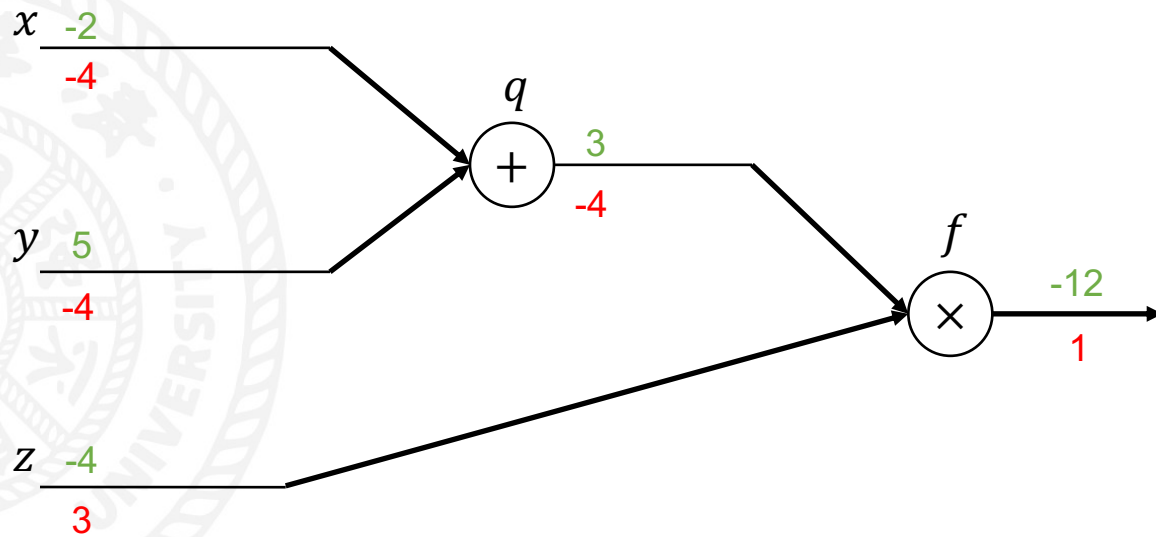


# Backpropagation



$$f(x, y, z) = (x + y)z$$

We want  $\frac{df}{dx}$ ,  $\frac{df}{dy}$ ,  $\frac{df}{dz}$

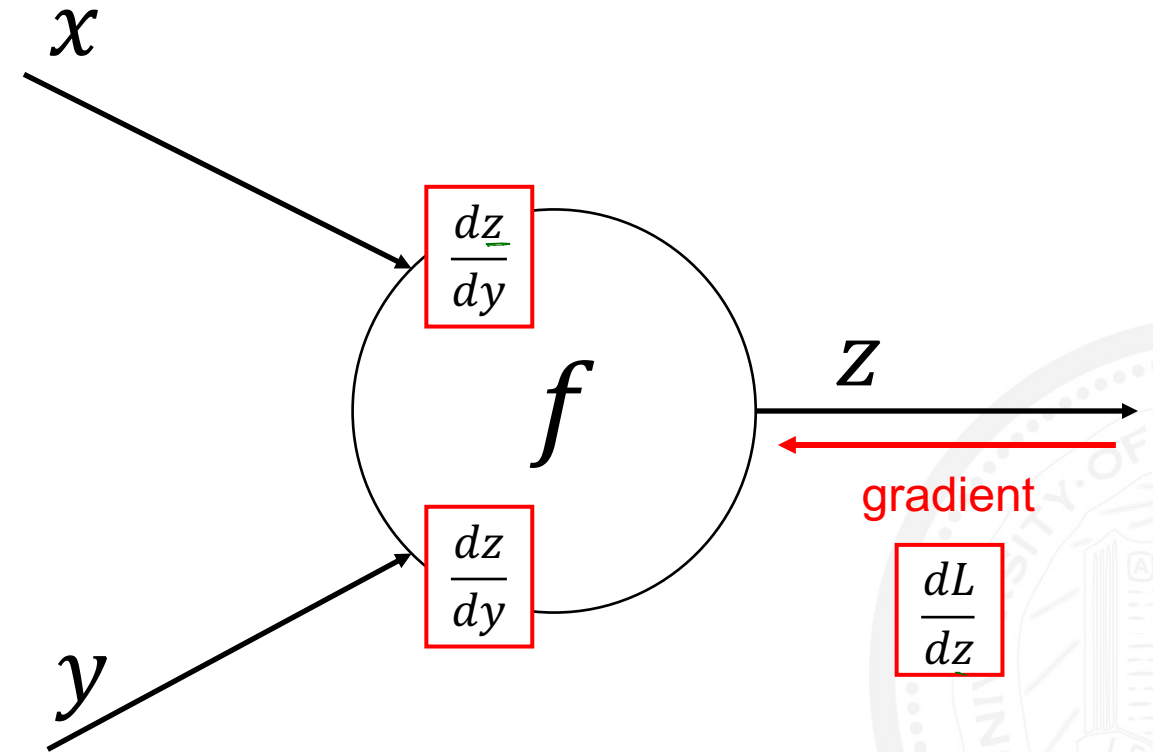
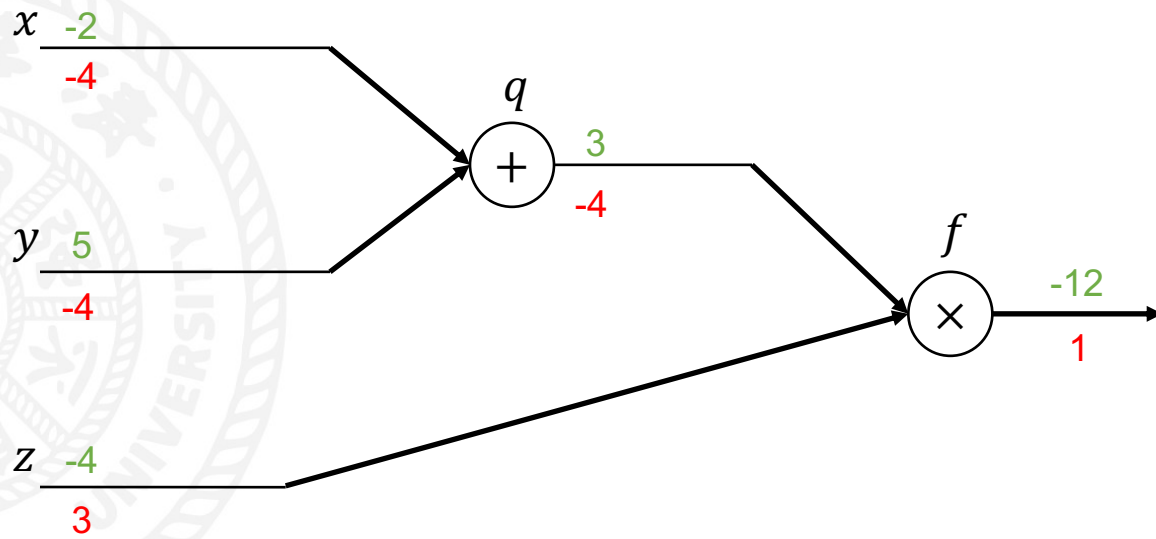


# Backpropagation



$$f(x, y, z) = (x + y)z$$

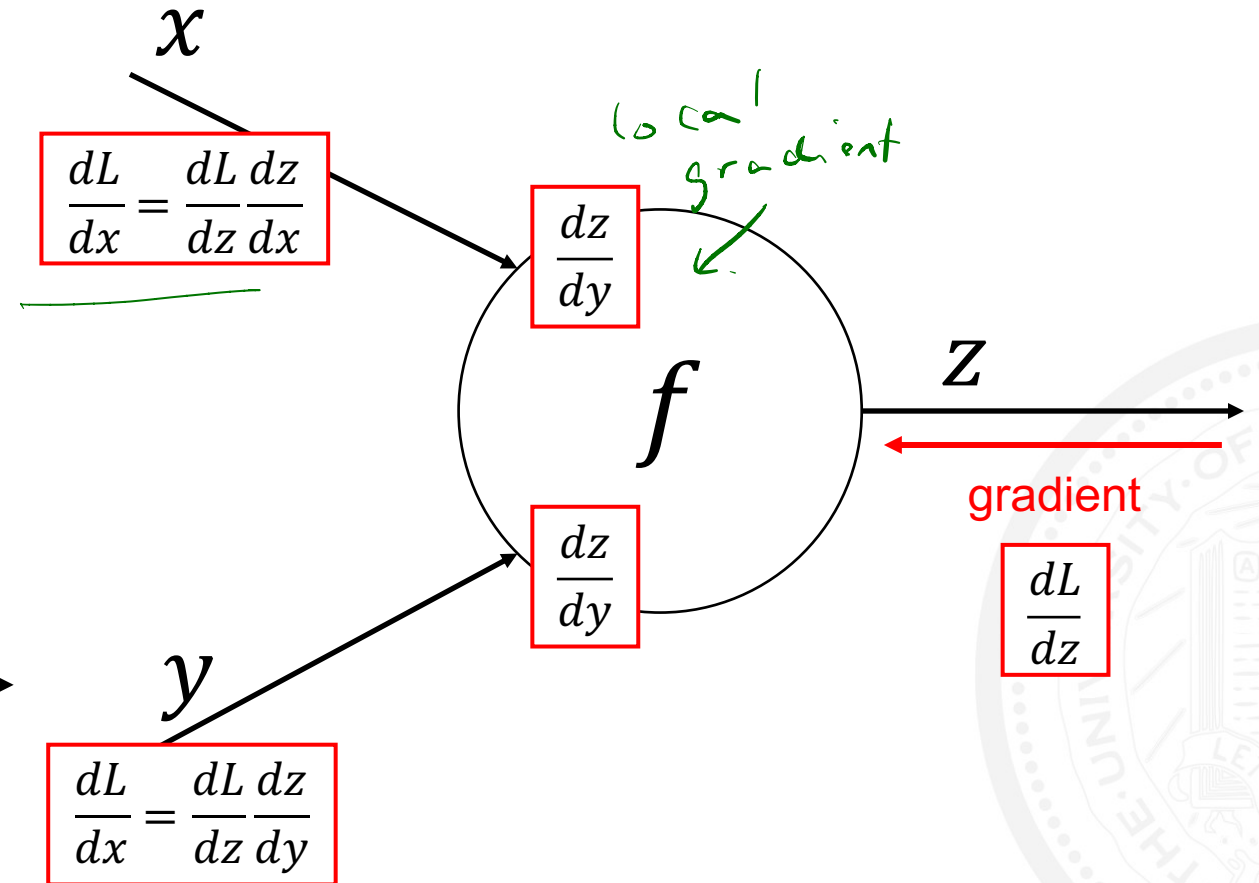
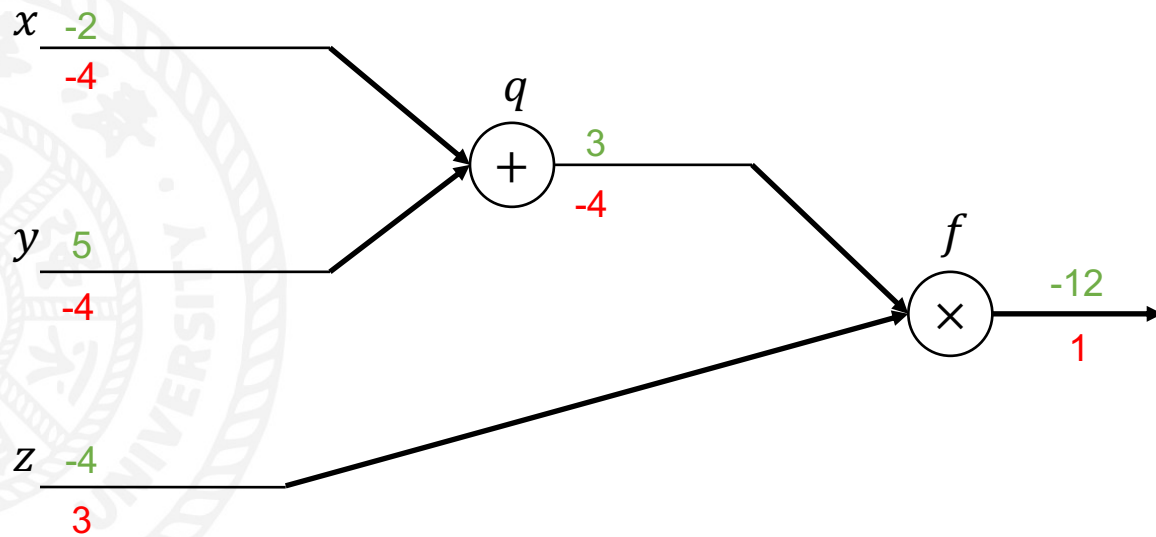
We want  $\frac{df}{dx}$ ,  $\frac{df}{dy}$ ,  $\frac{df}{dz}$



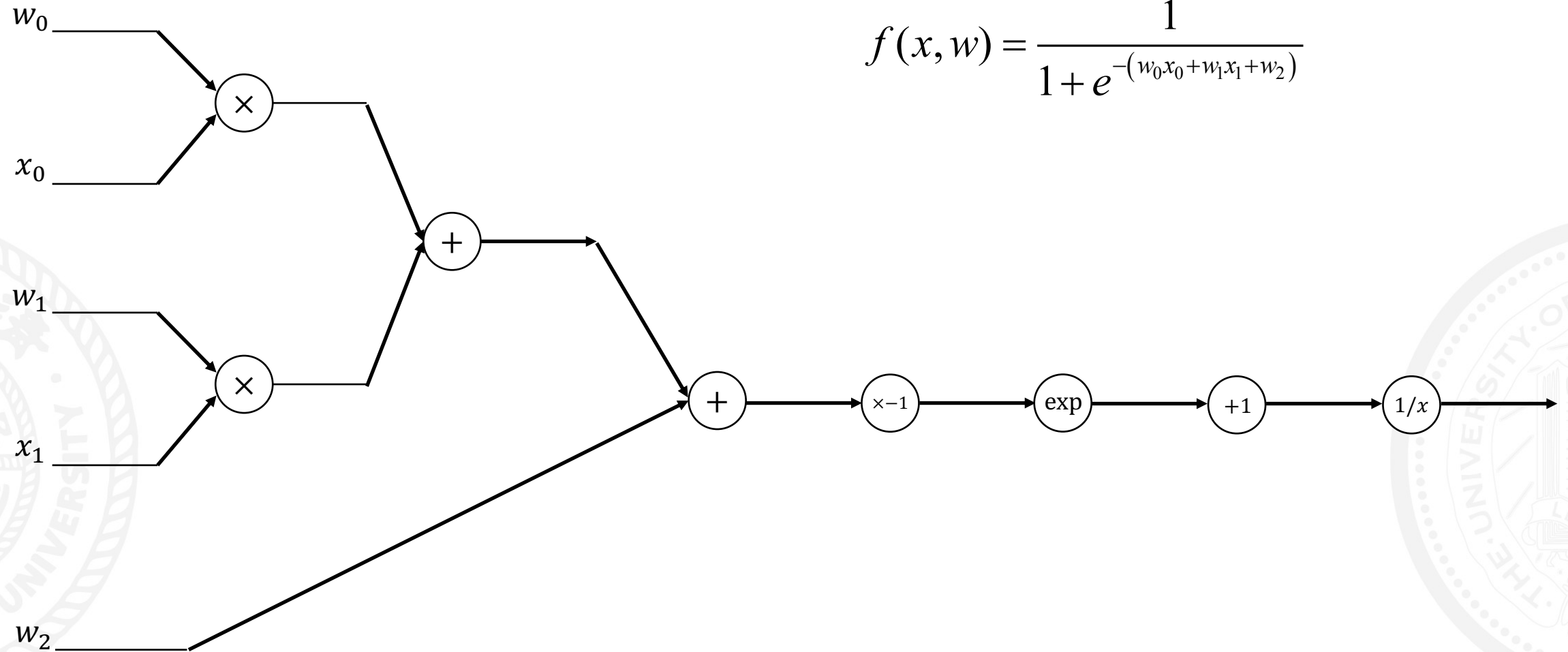
# Backpropagation

$$f(x, y, z) = (x + y)z$$

We want  $\frac{df}{dx}$ ,  $\frac{df}{dy}$ ,  $\frac{df}{dz}$

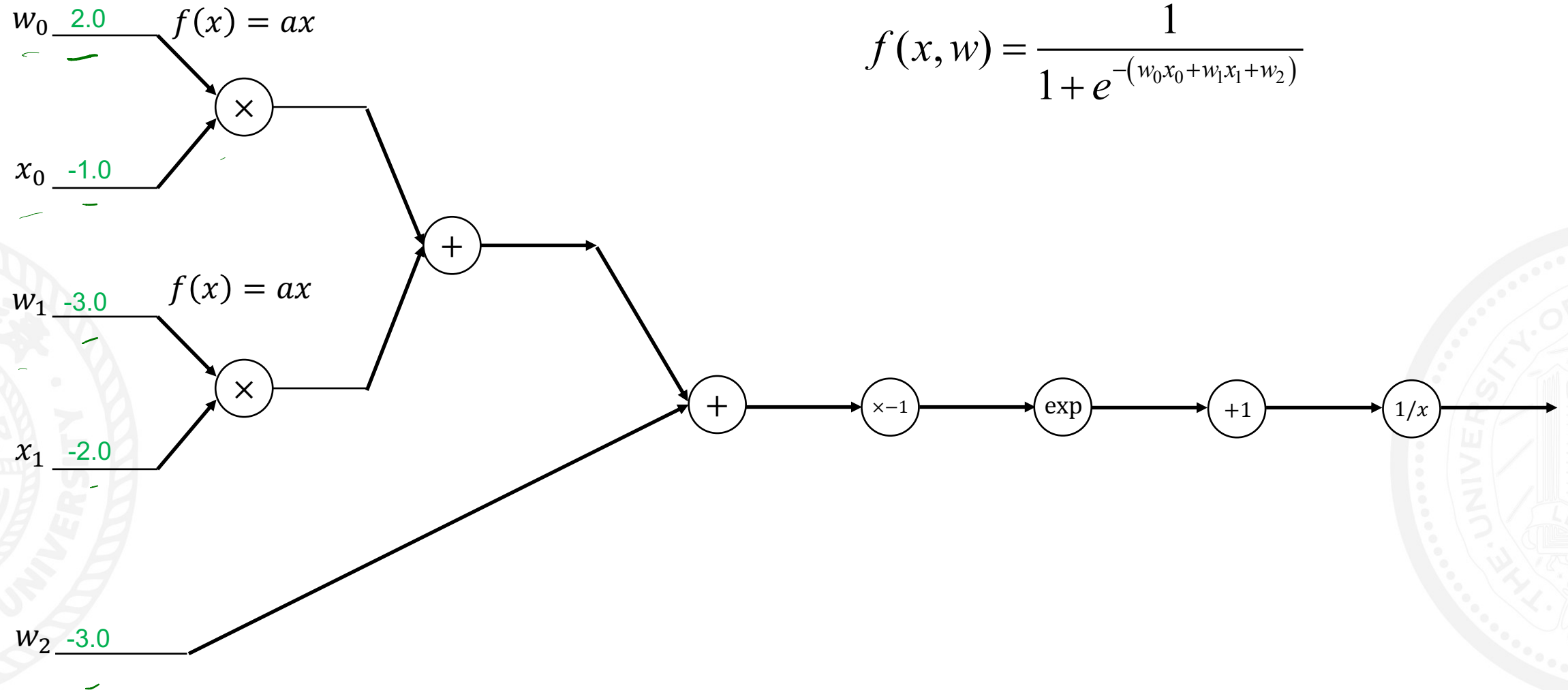


# Backpropagation



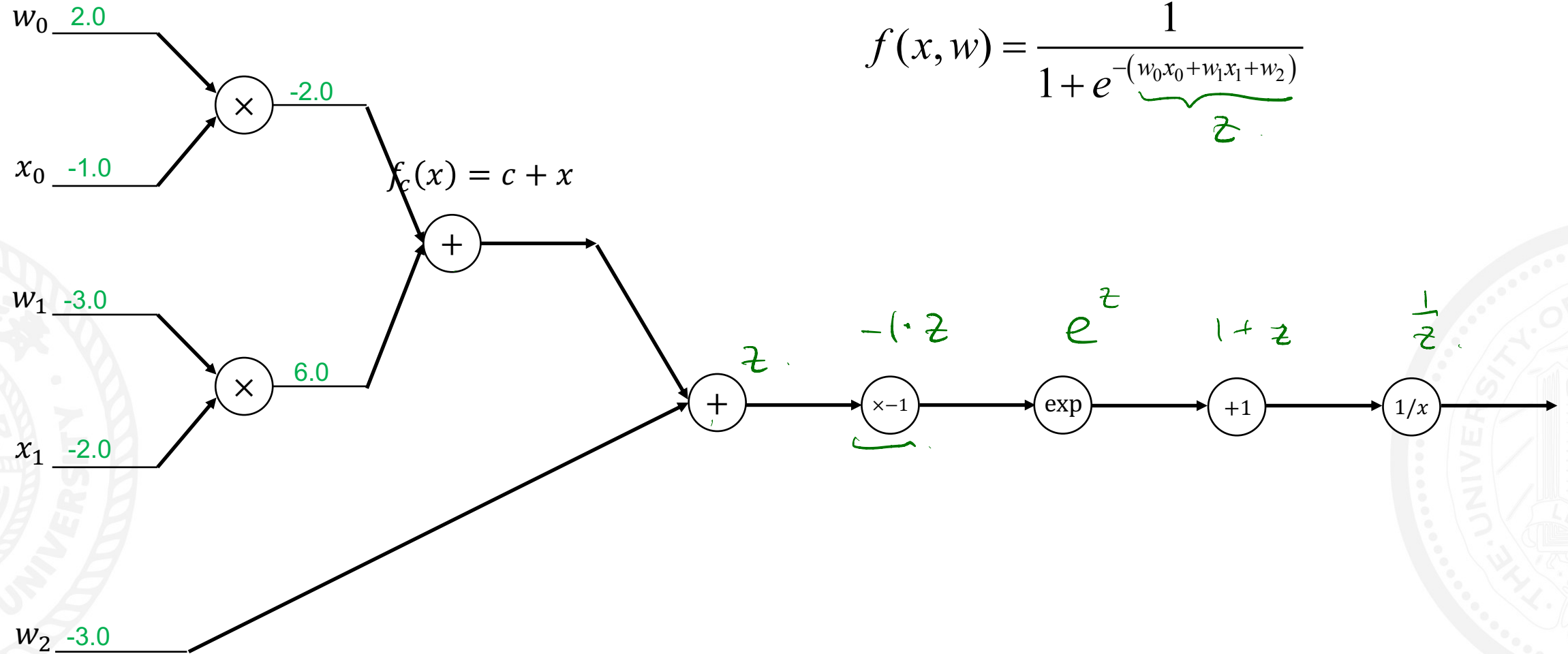
$$f(x, w) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

# Backpropagation

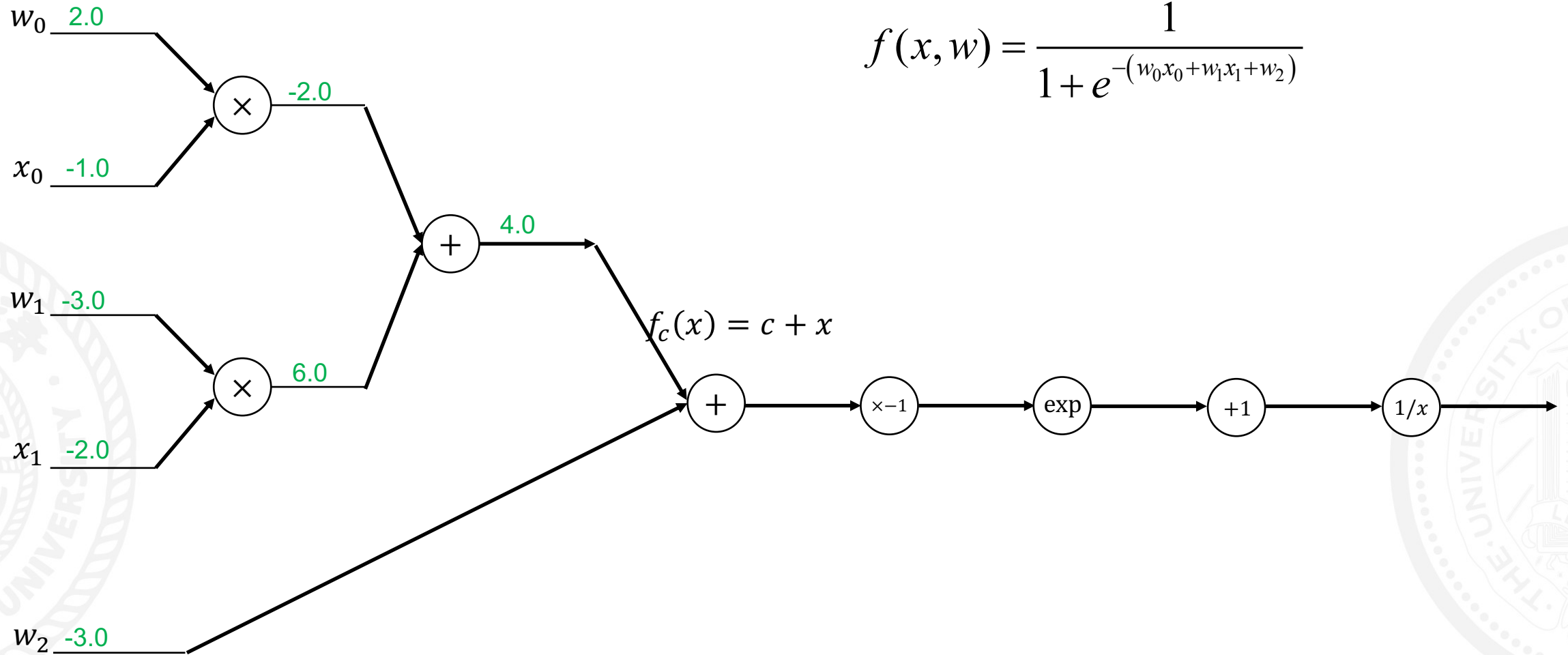


$$f(x, w) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

# Backpropagation



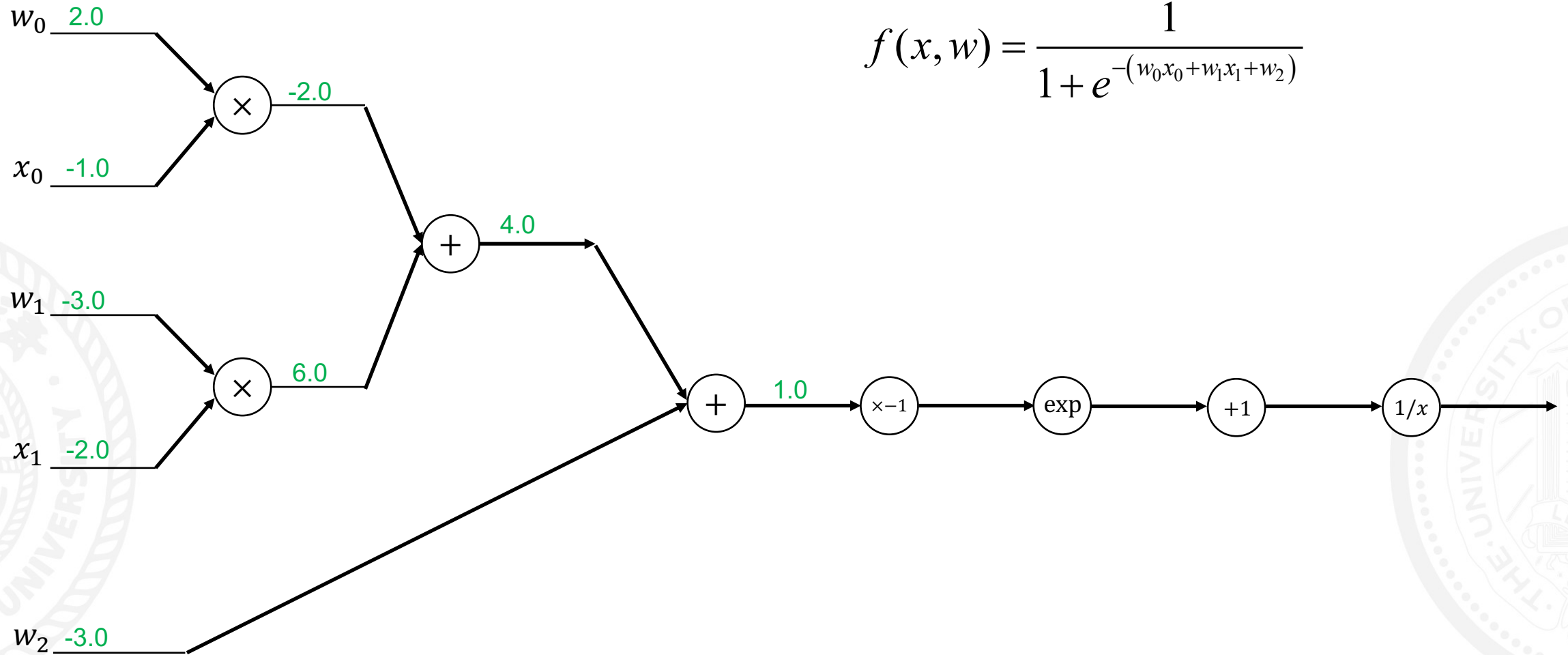
# Backpropagation



$$f(x, w) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

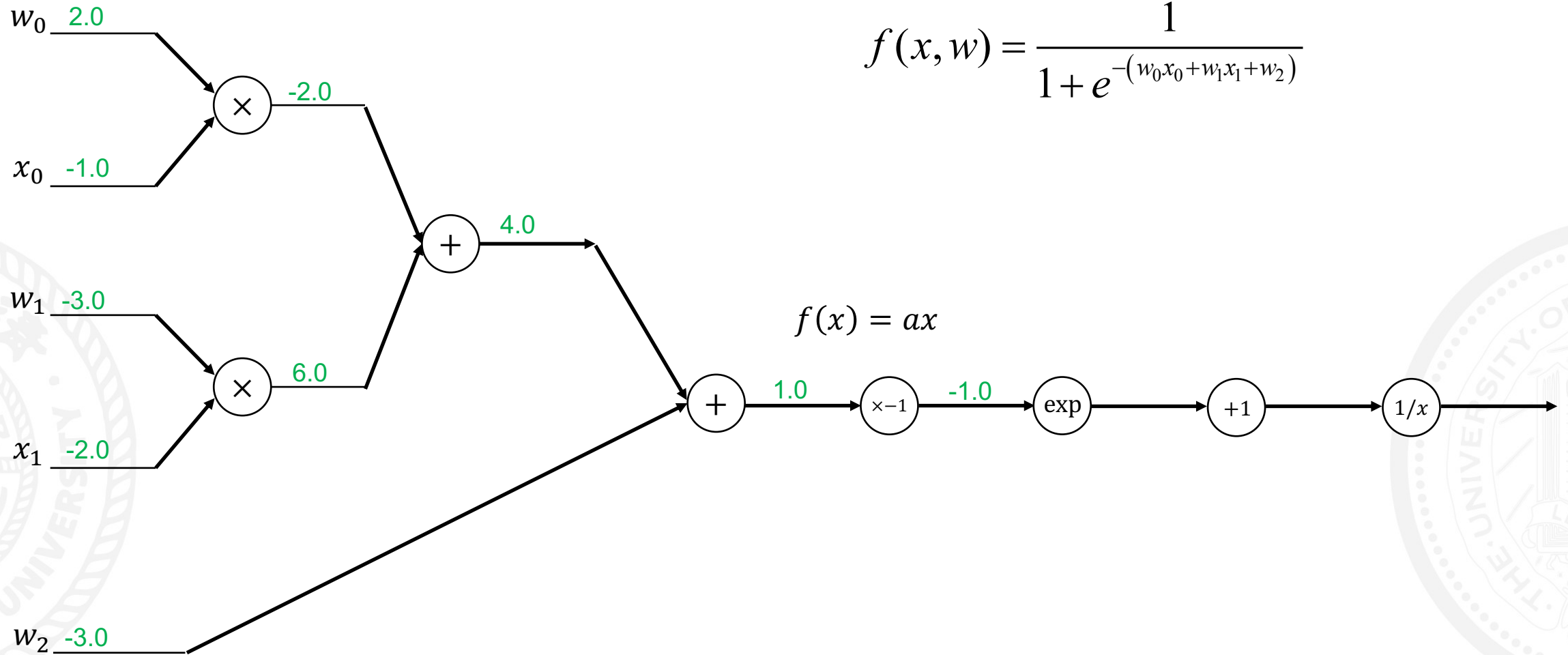


# Backpropagation



$$f(x, w) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

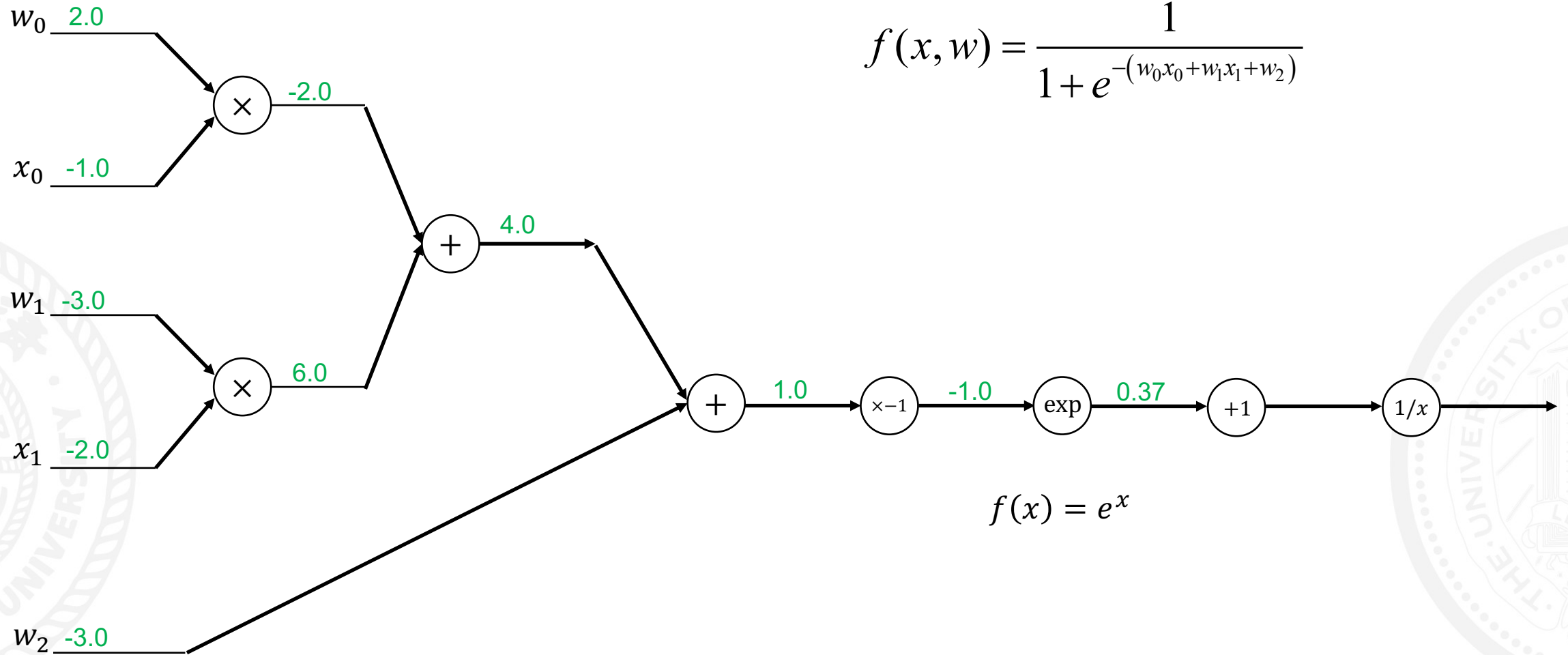
# Backpropagation



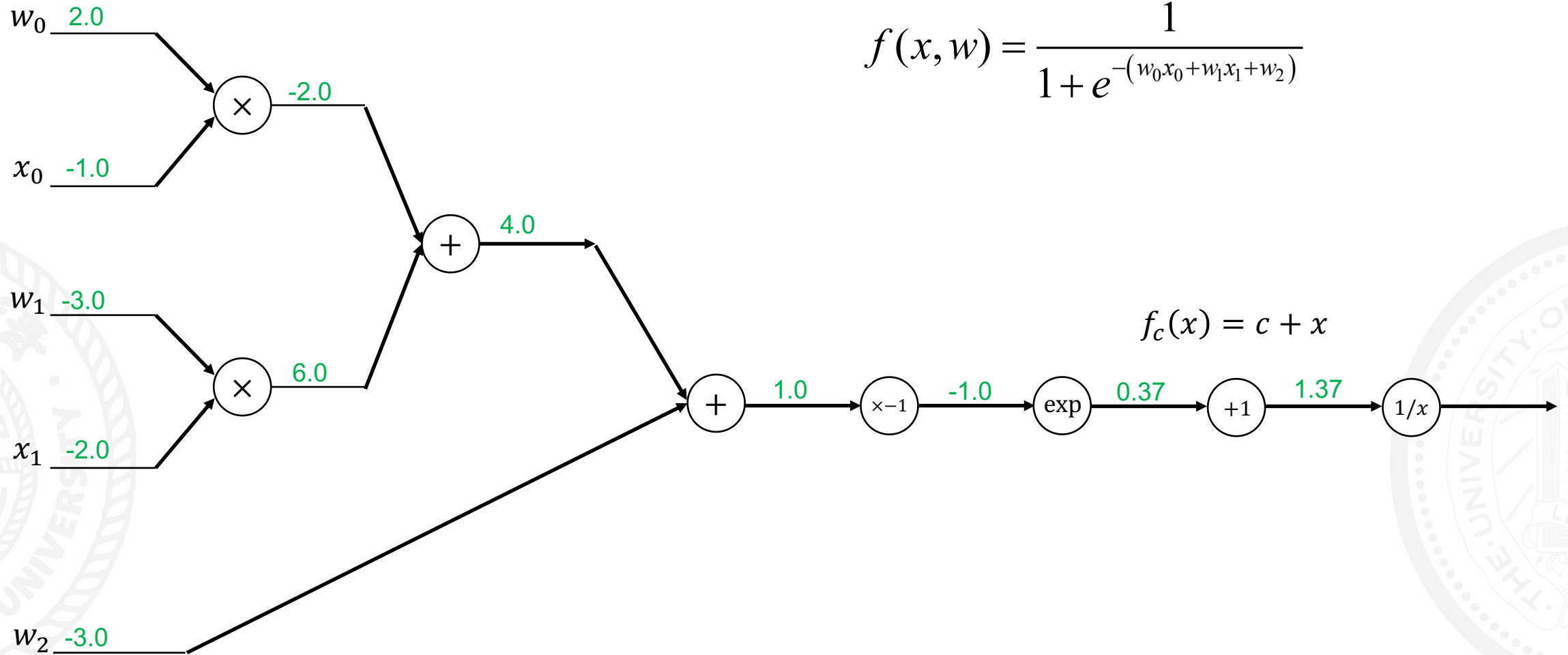
$$f(x, w) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

$$f(x) = ax$$

# Backpropagation



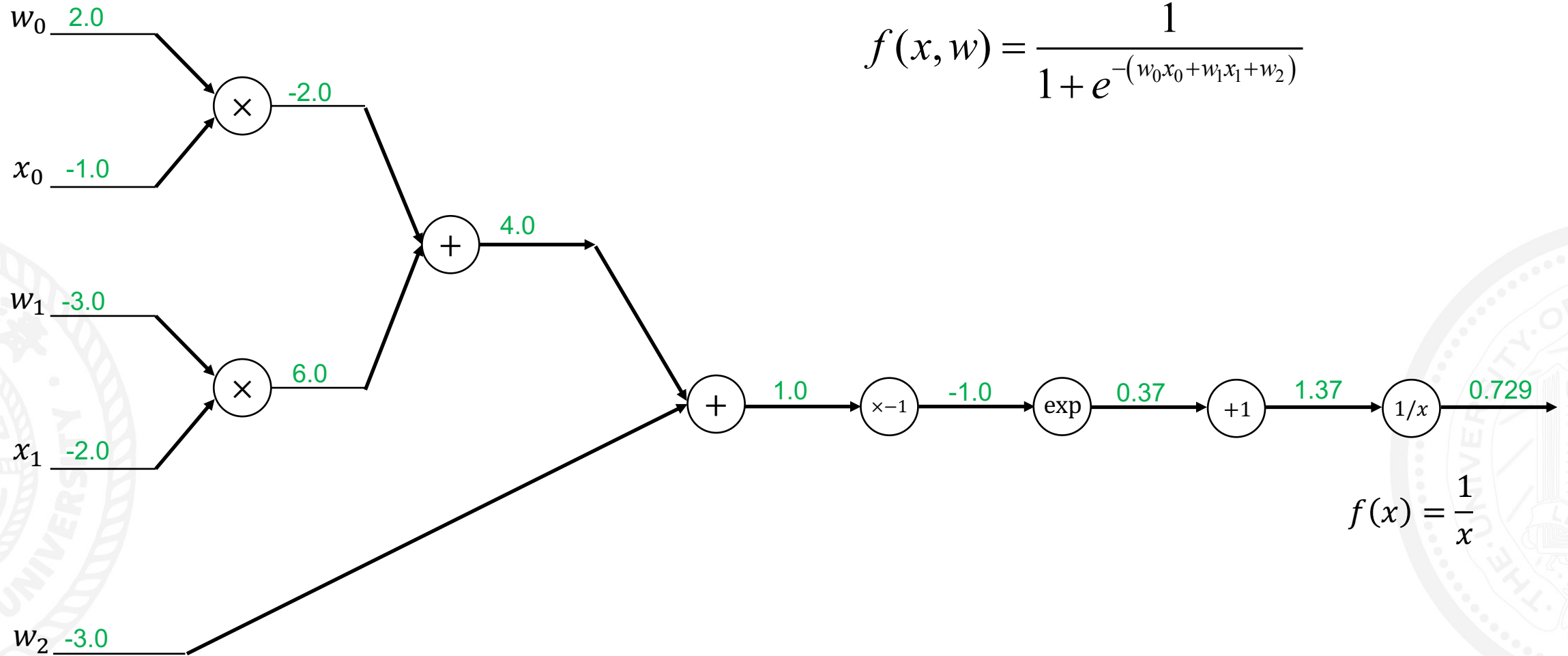
# Backpropagation



$$f(x, w) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

$$f_c(x) = c + x$$

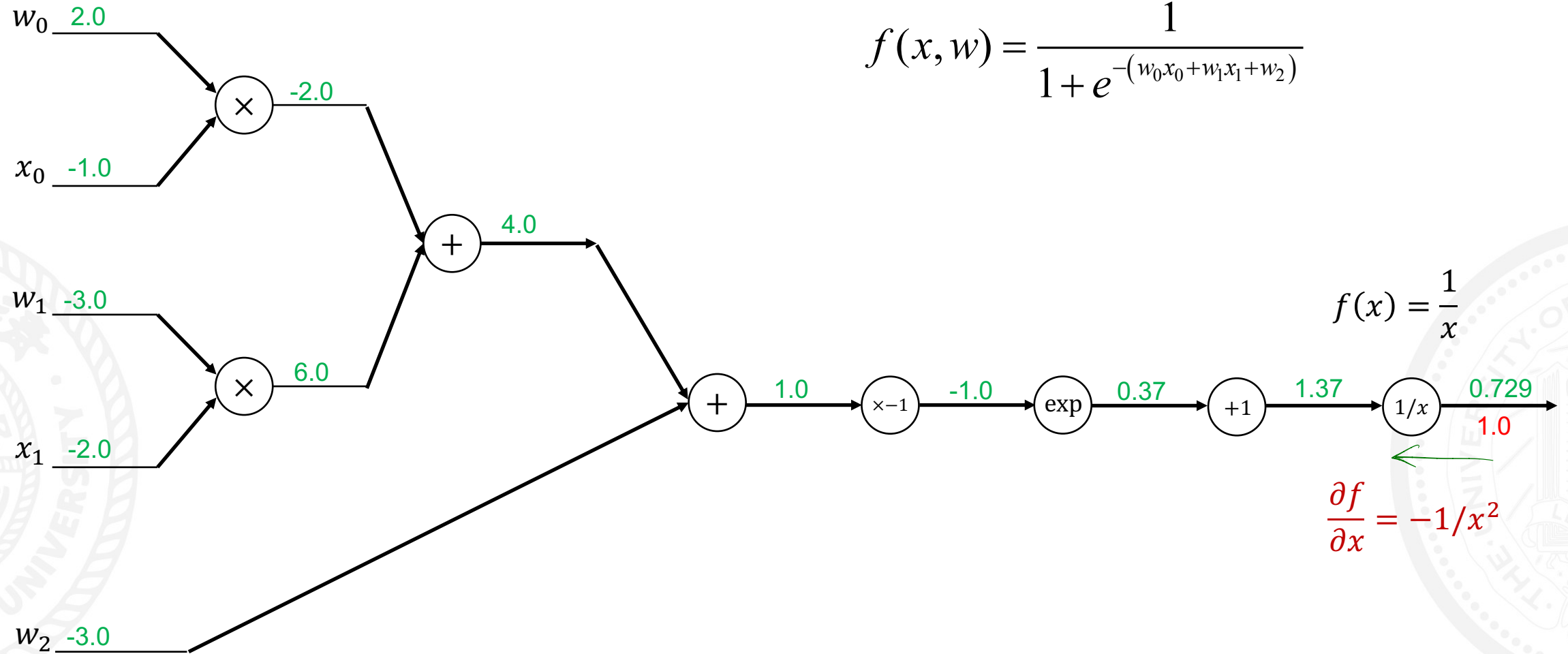
# Backpropagation



$$f(x, w) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

$$f(x) = \frac{1}{x}$$

# Backpropagation

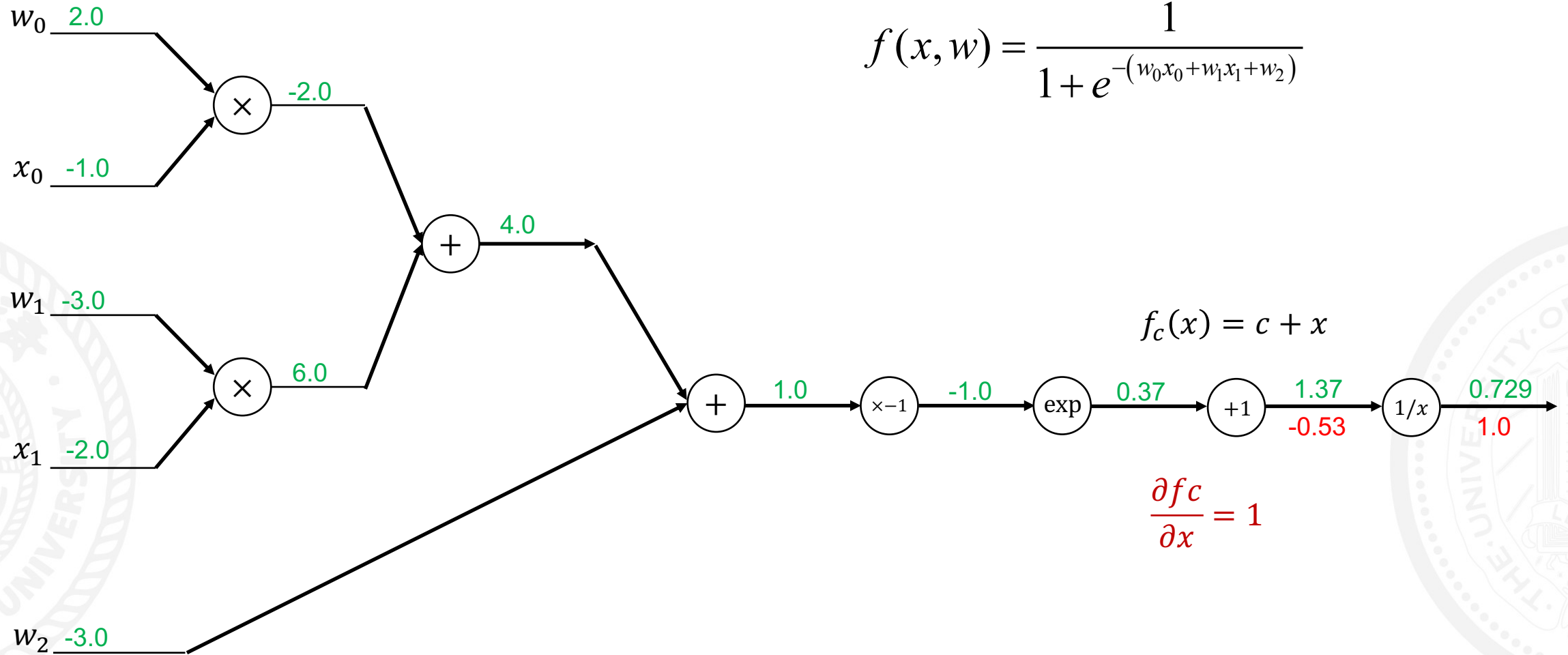


$$f(x, w) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

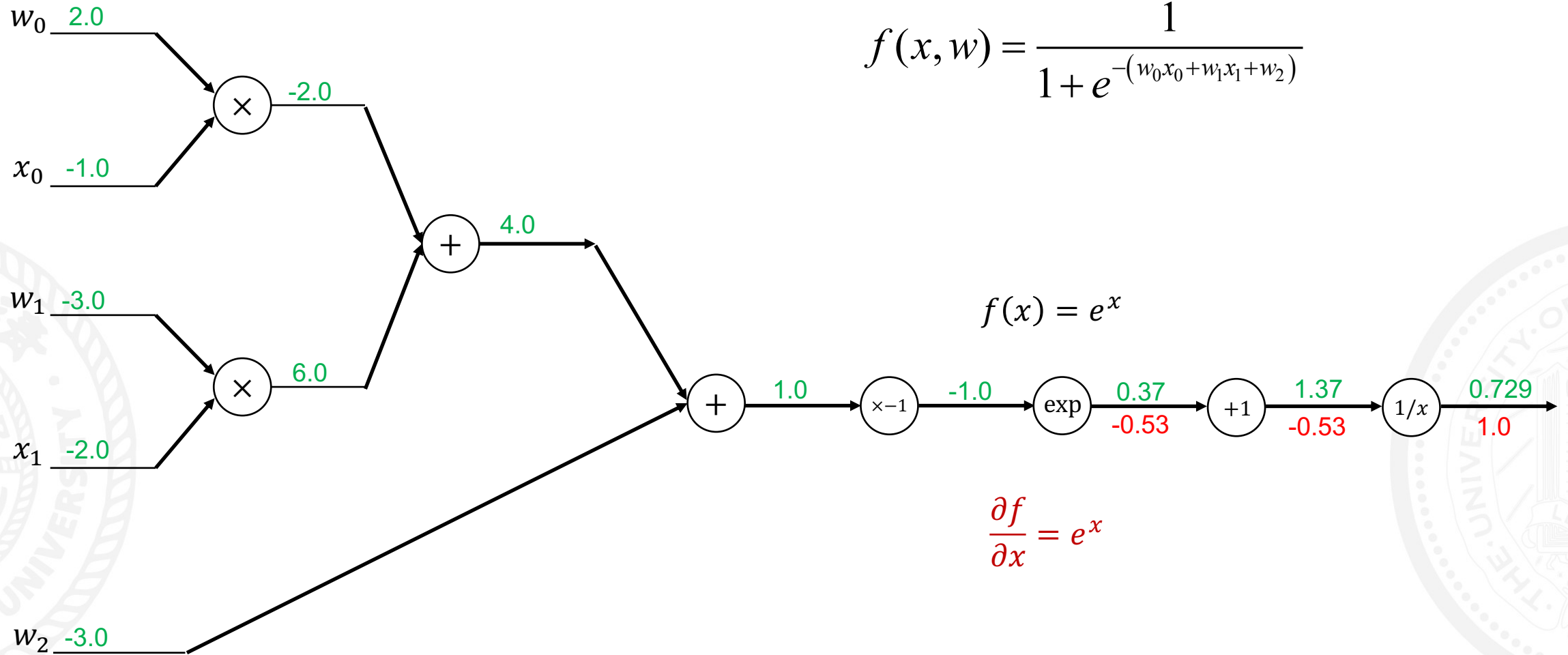
$$f(x) = \frac{1}{x}$$

$$\frac{\partial f}{\partial x} = -\frac{1}{x^2}$$

# Backpropagation

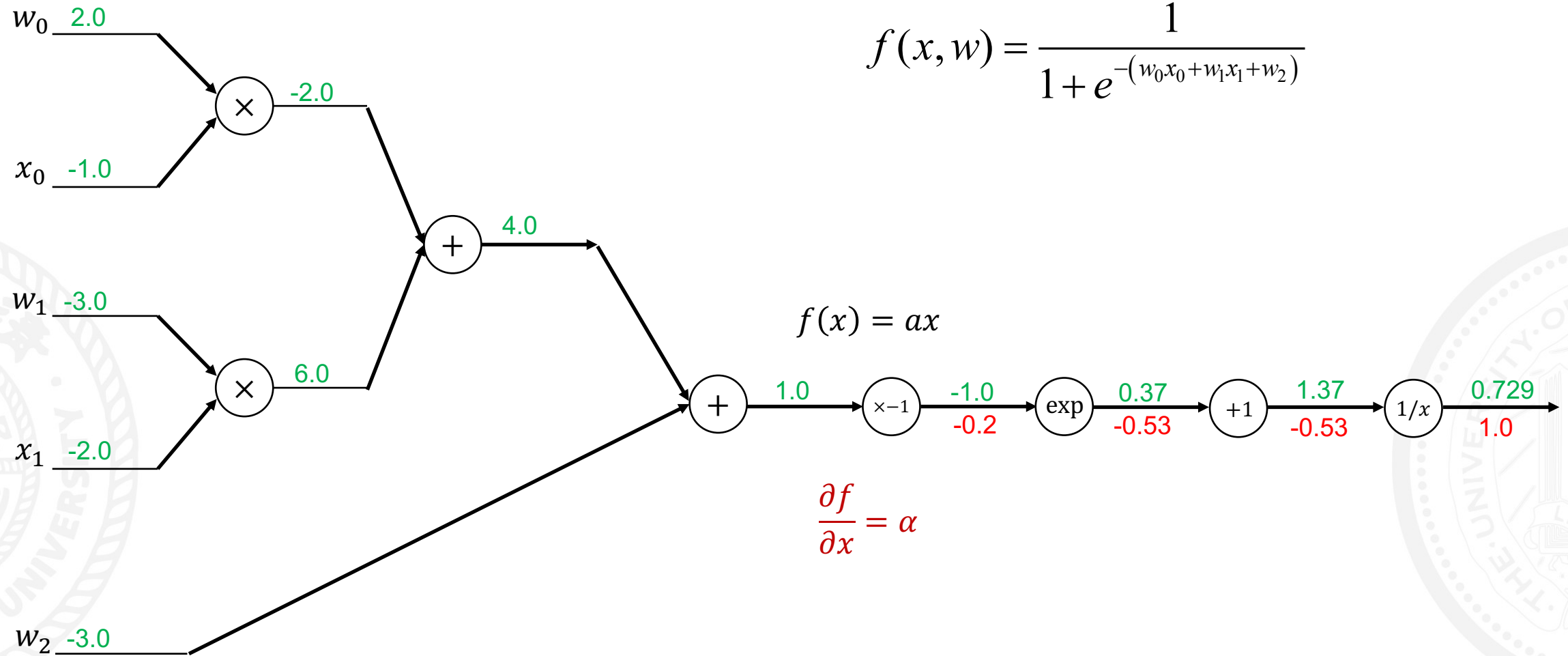


# Backpropagation

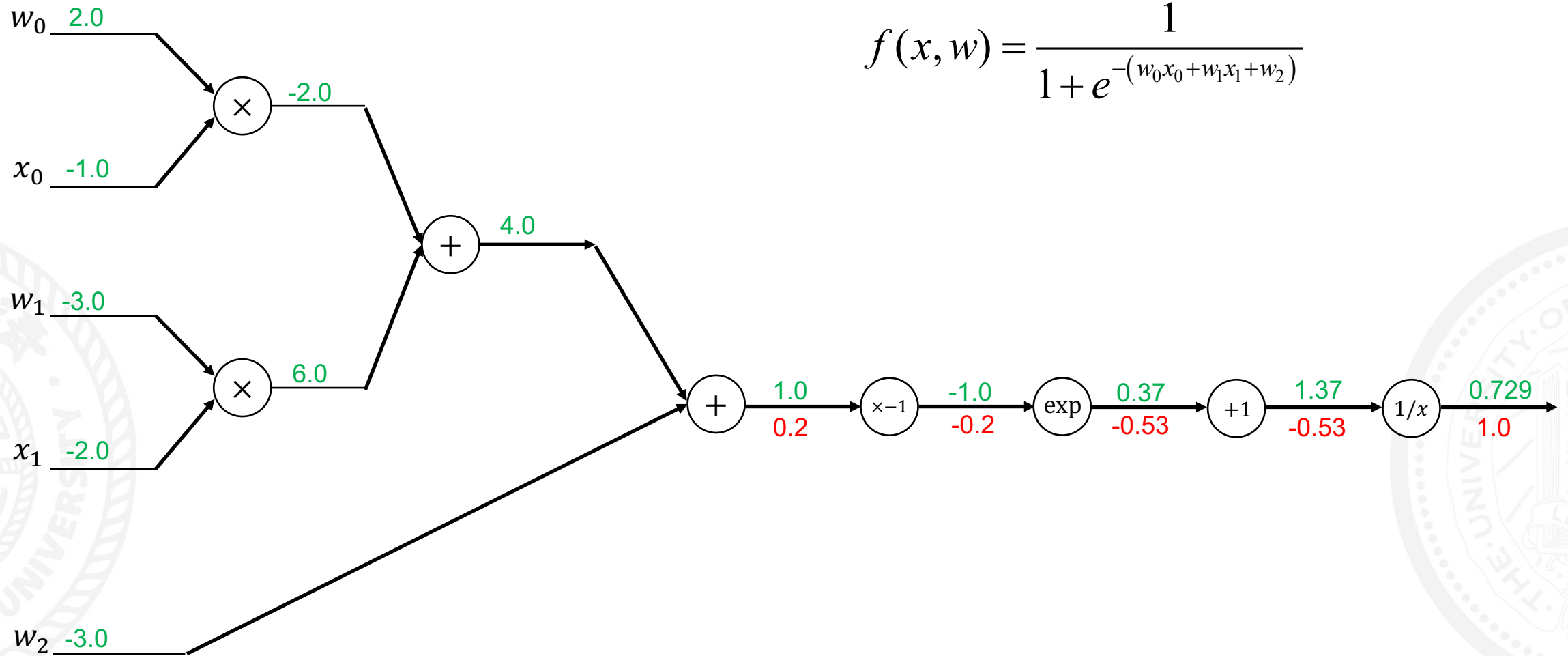




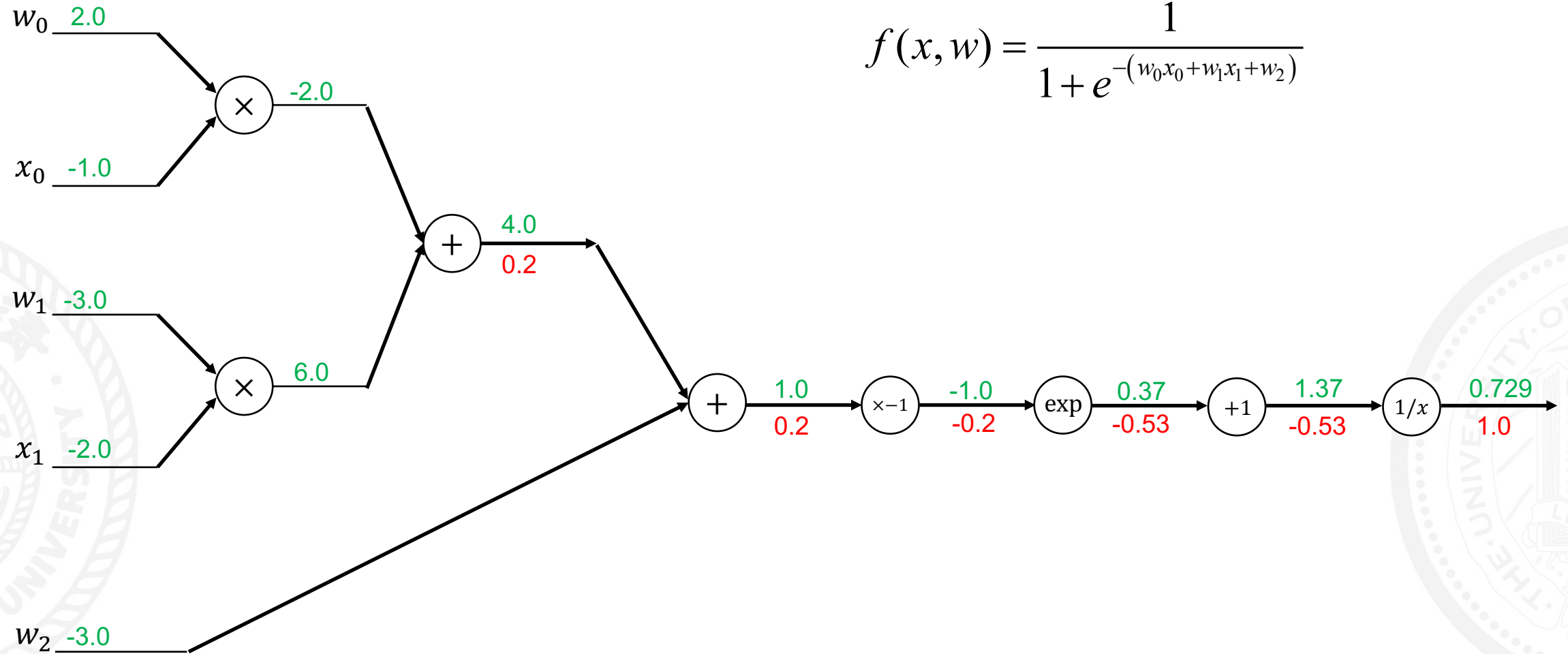
# Backpropagation



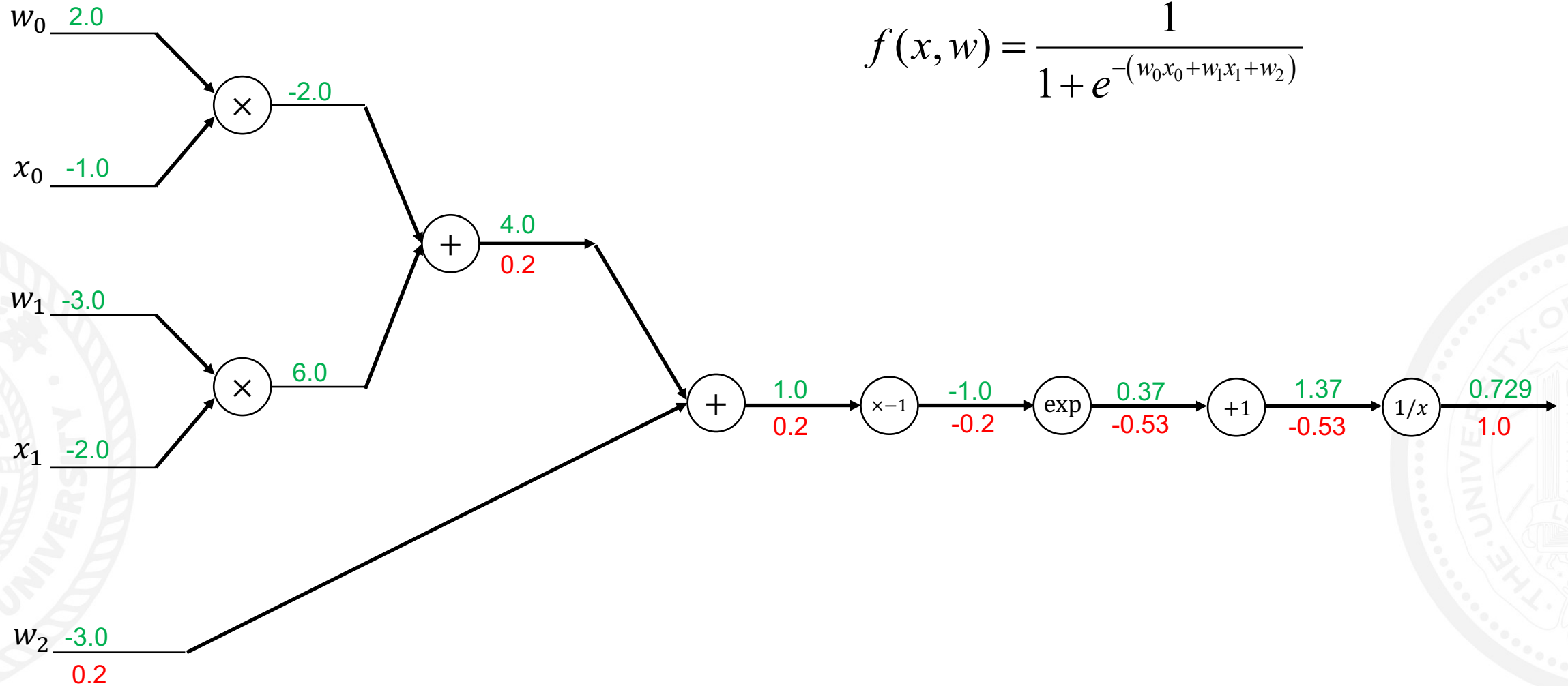
# Backpropagation



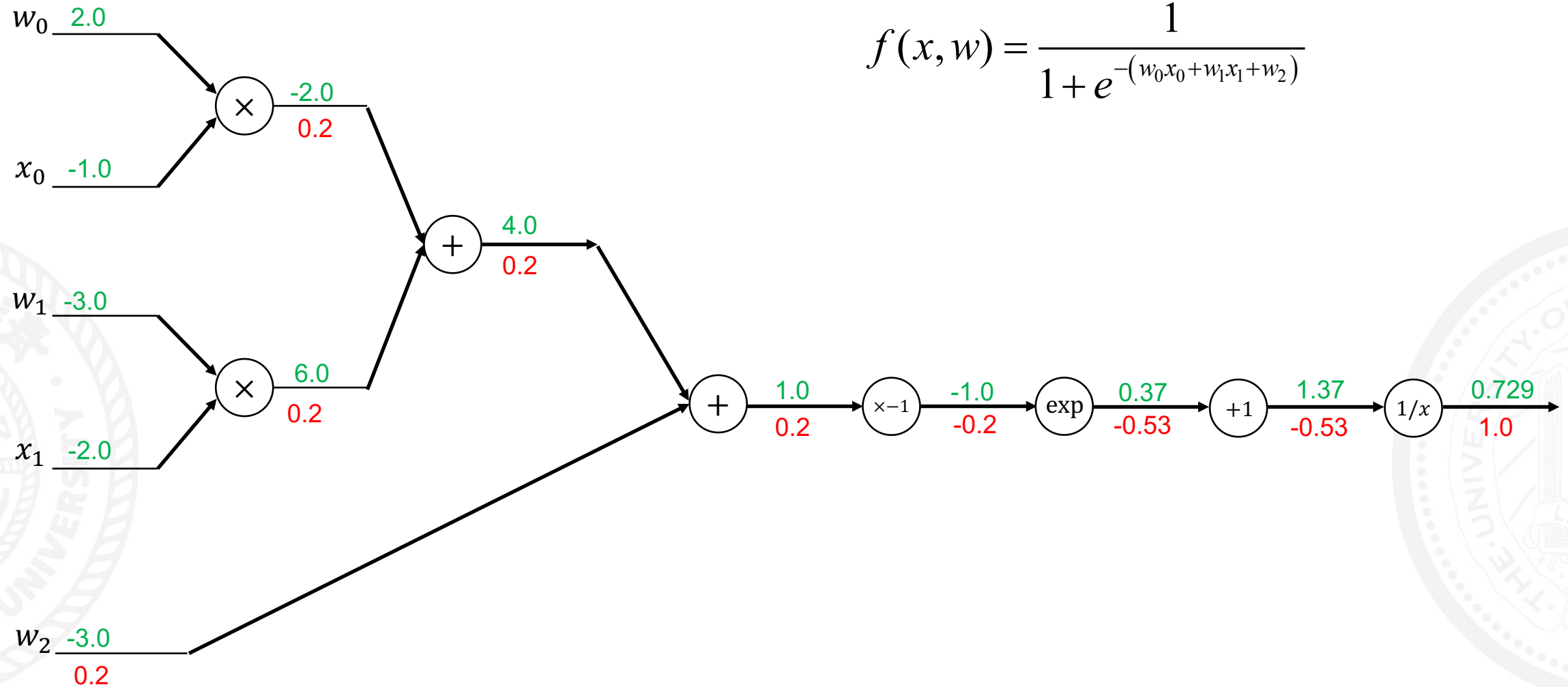
# Backpropagation



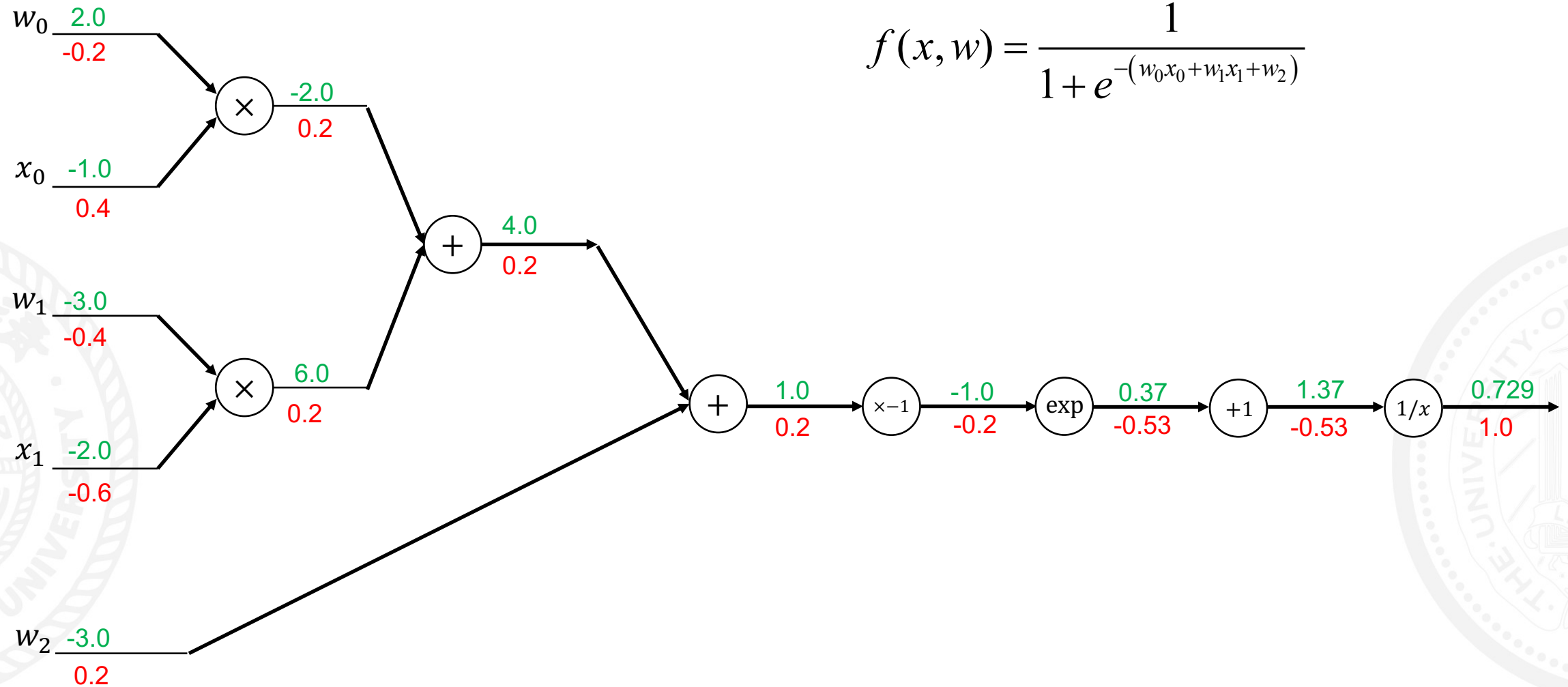
# Backpropagation



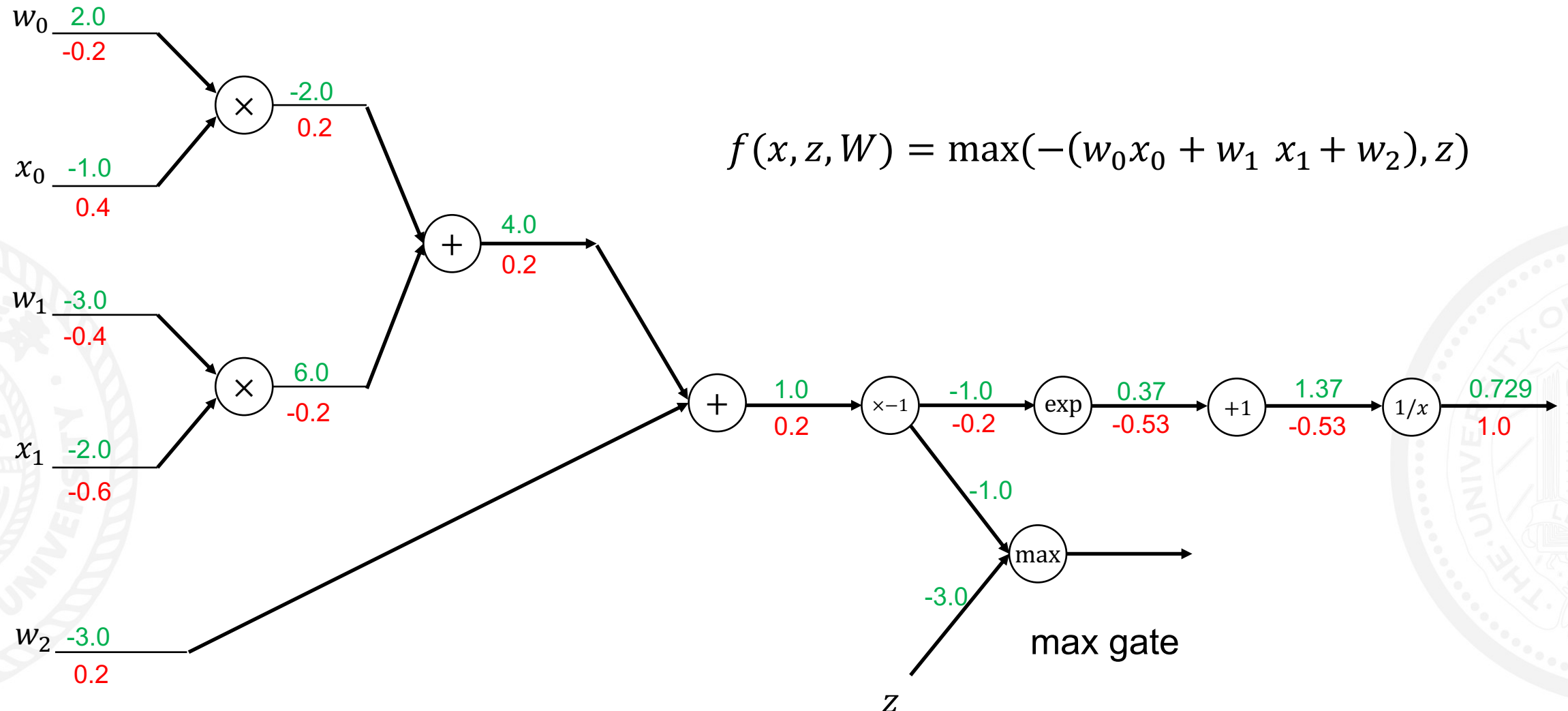
# Backpropagation



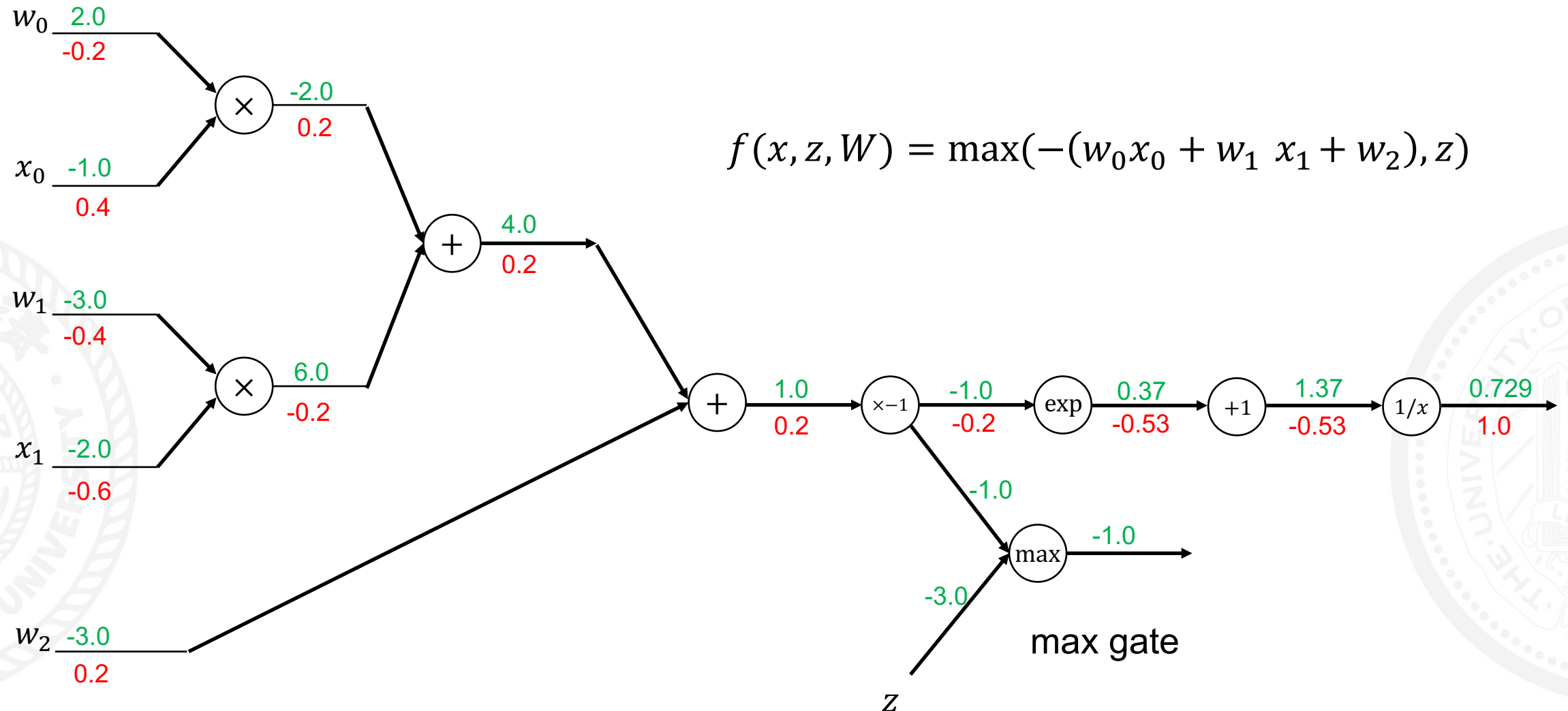
# Backpropagation



# Backpropagation

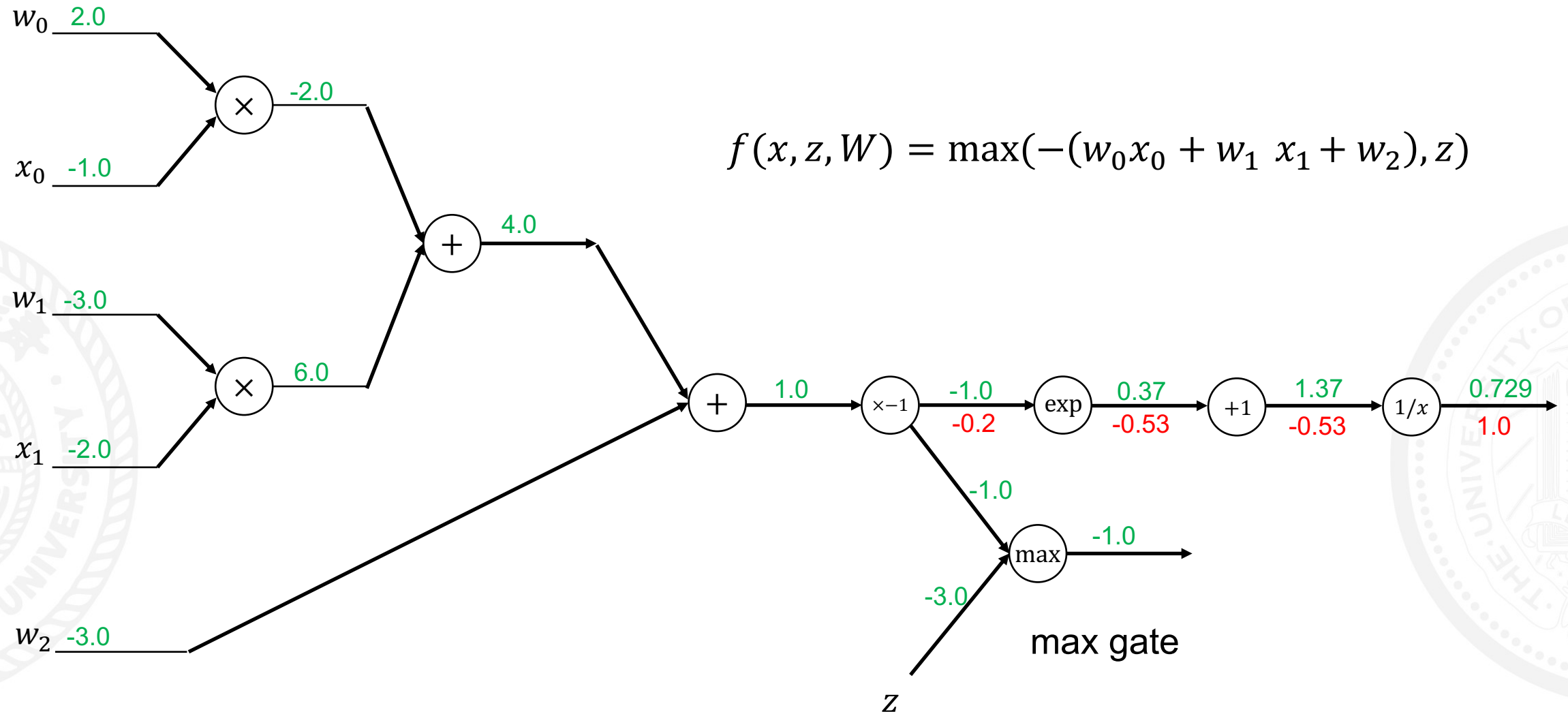


# Backpropagation

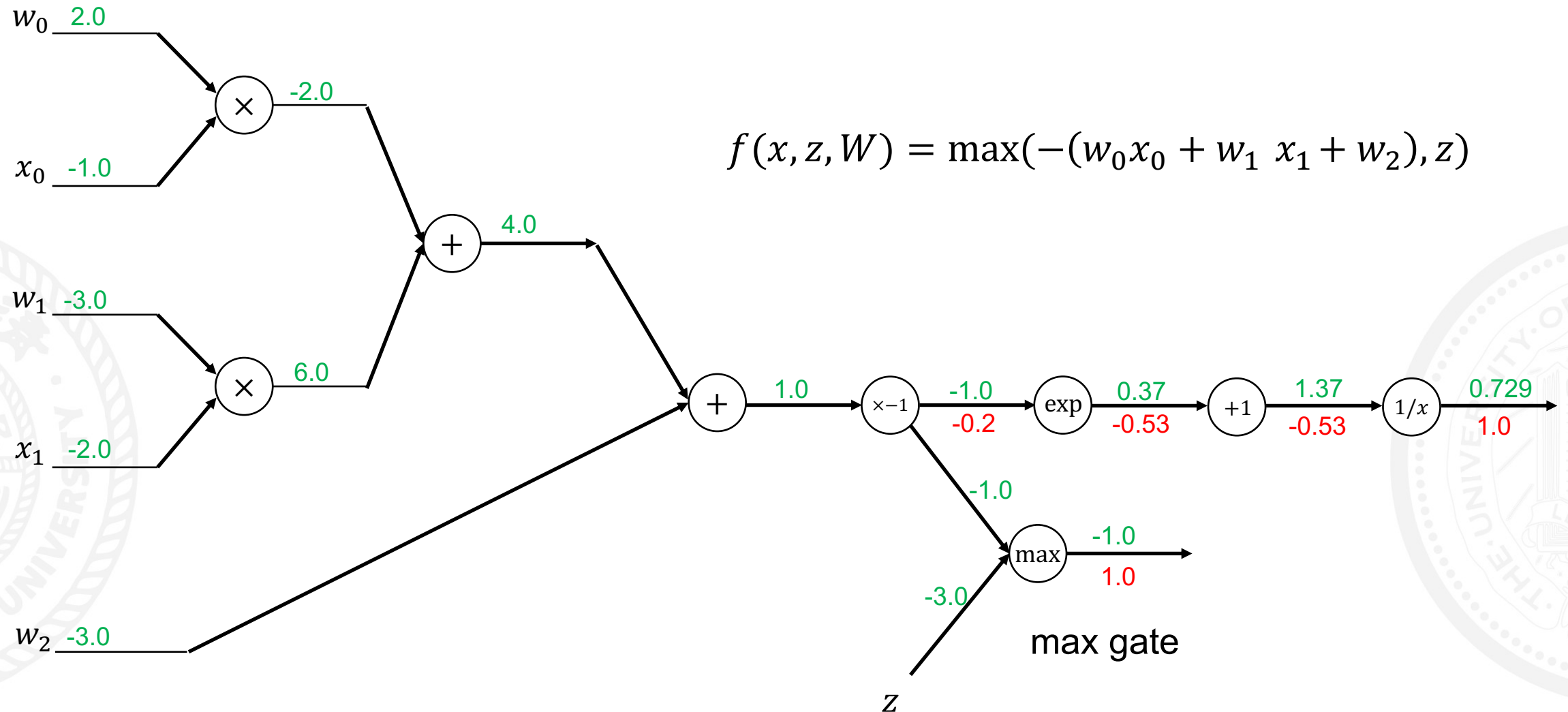




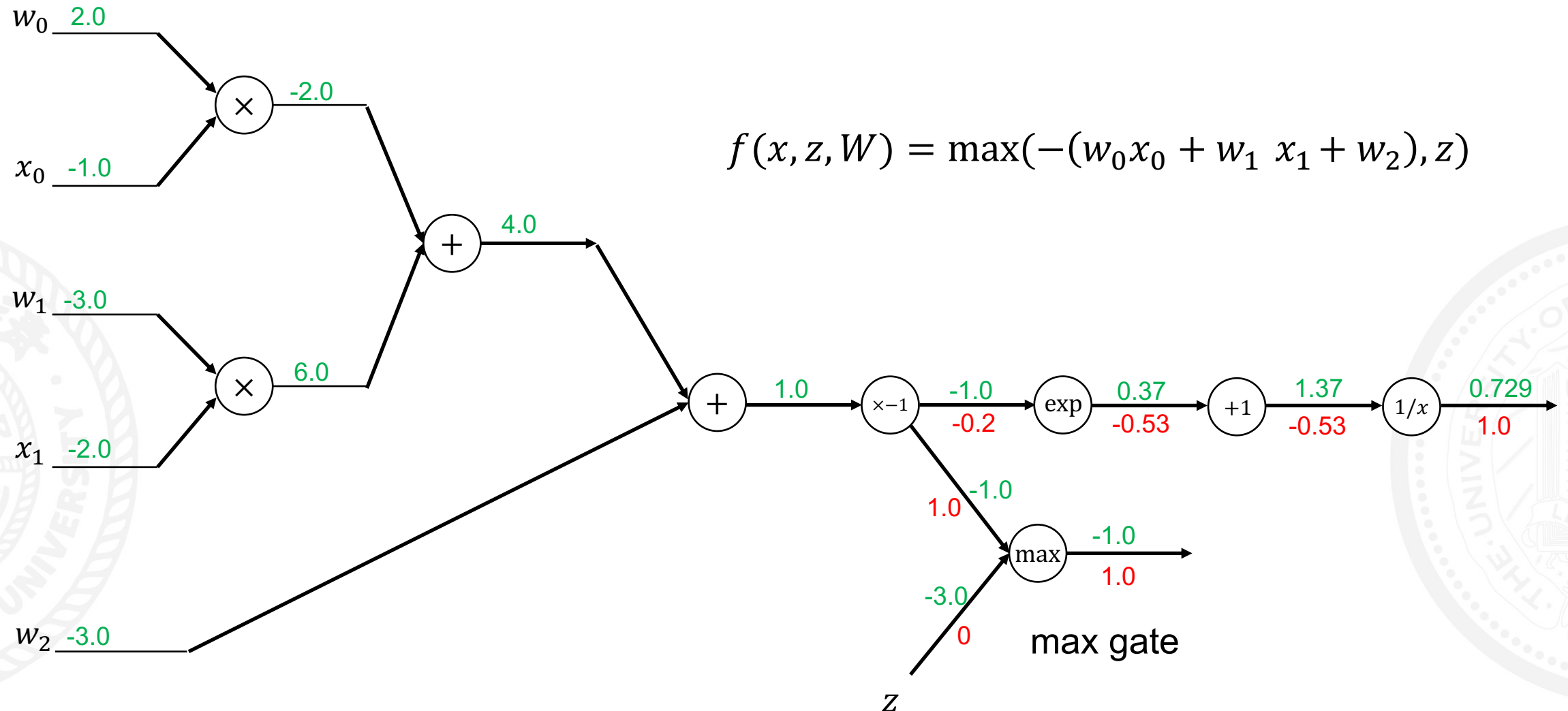
# Backpropagation



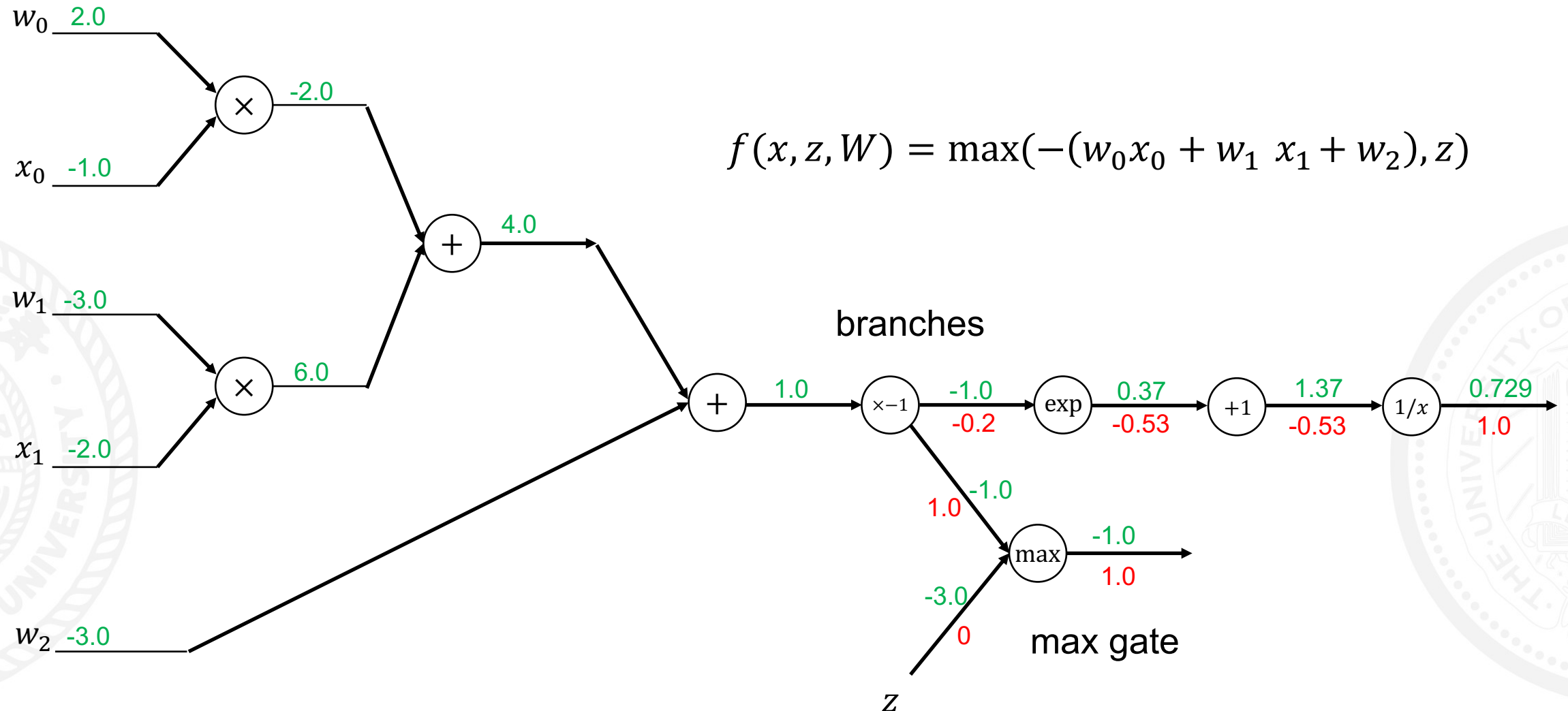
# Backpropagation



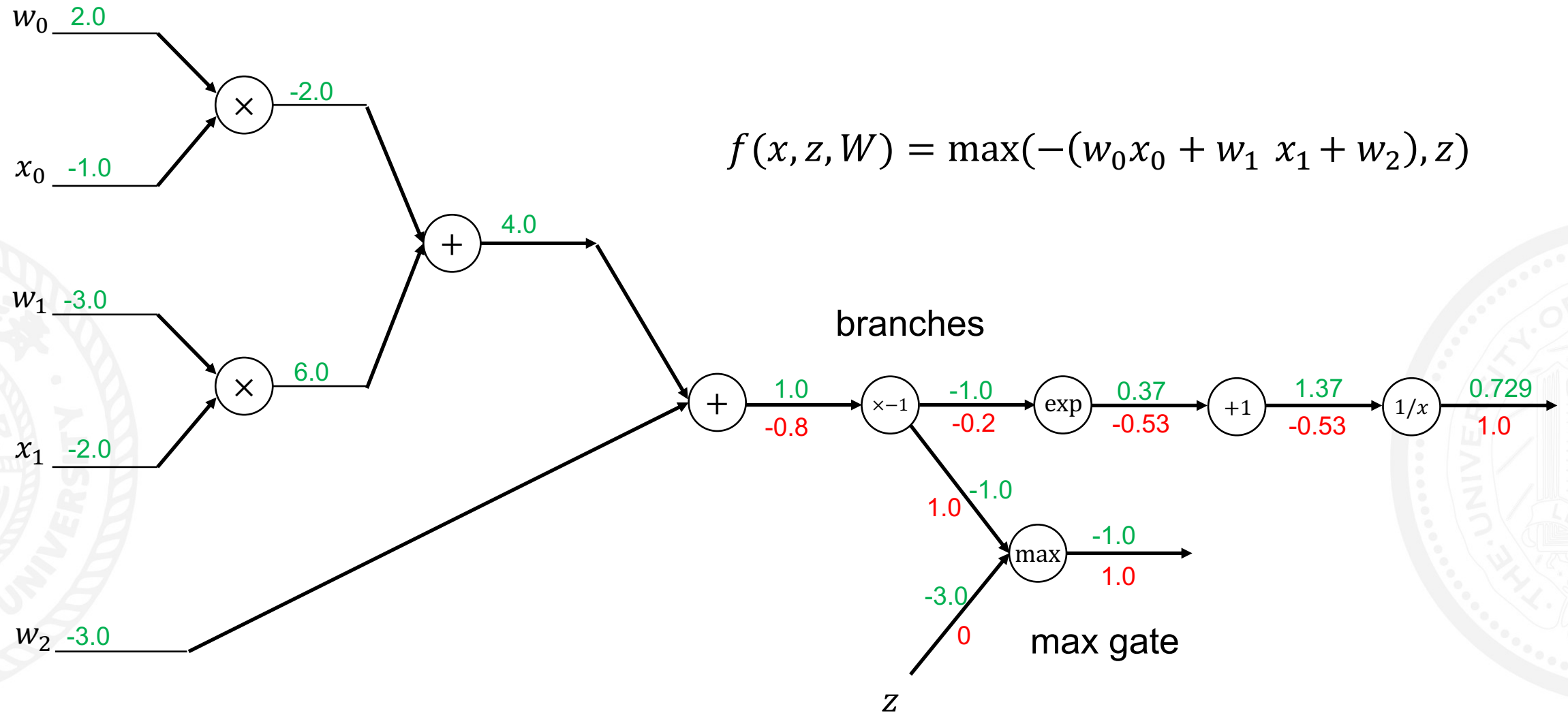
# Backpropagation



# Backpropagation



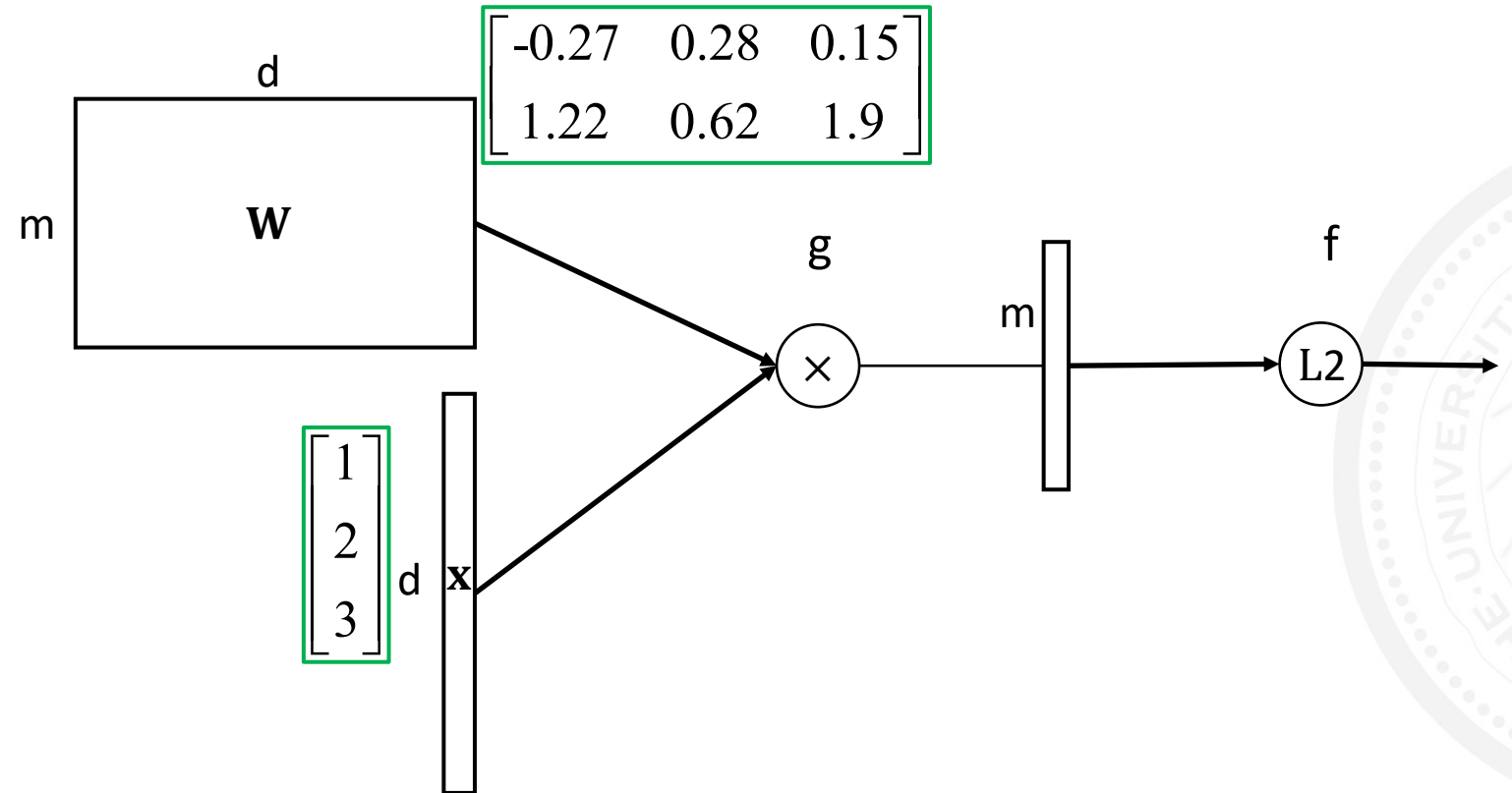
# Backpropagation



# Backpropagation

Vectorized example

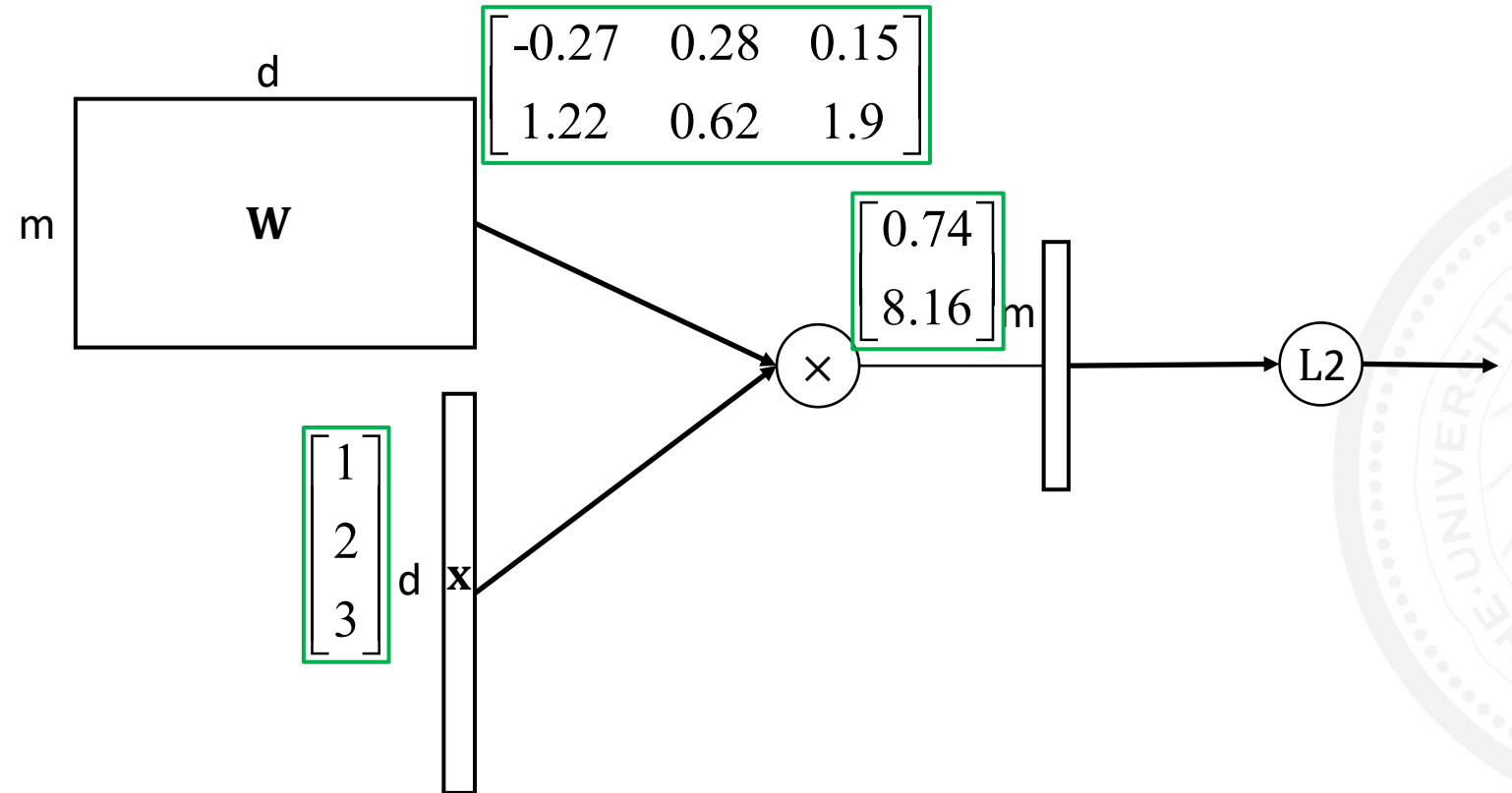
$$f(\mathbf{x}, \mathbf{W}) = \|\mathbf{W} \cdot \mathbf{x}\|^2$$



# Backpropagation

Vectorized example

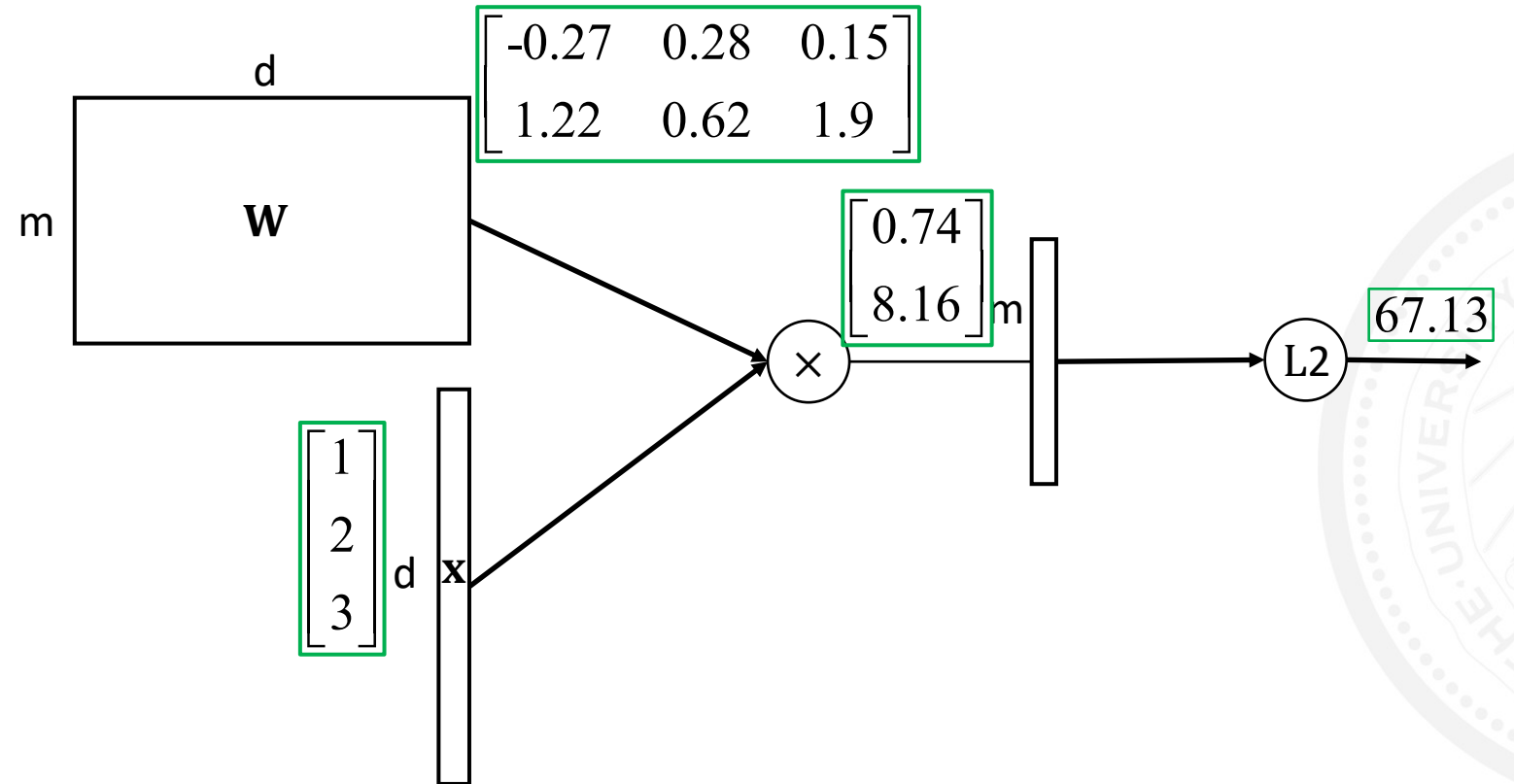
$$f(\mathbf{x}, \mathbf{W}) = \|\mathbf{W} \cdot \mathbf{x}\|^2$$



# Backpropagation

Vectorized example

$$f(\mathbf{x}, \mathbf{W}) = \|\mathbf{W} \cdot \mathbf{x}\|^2$$

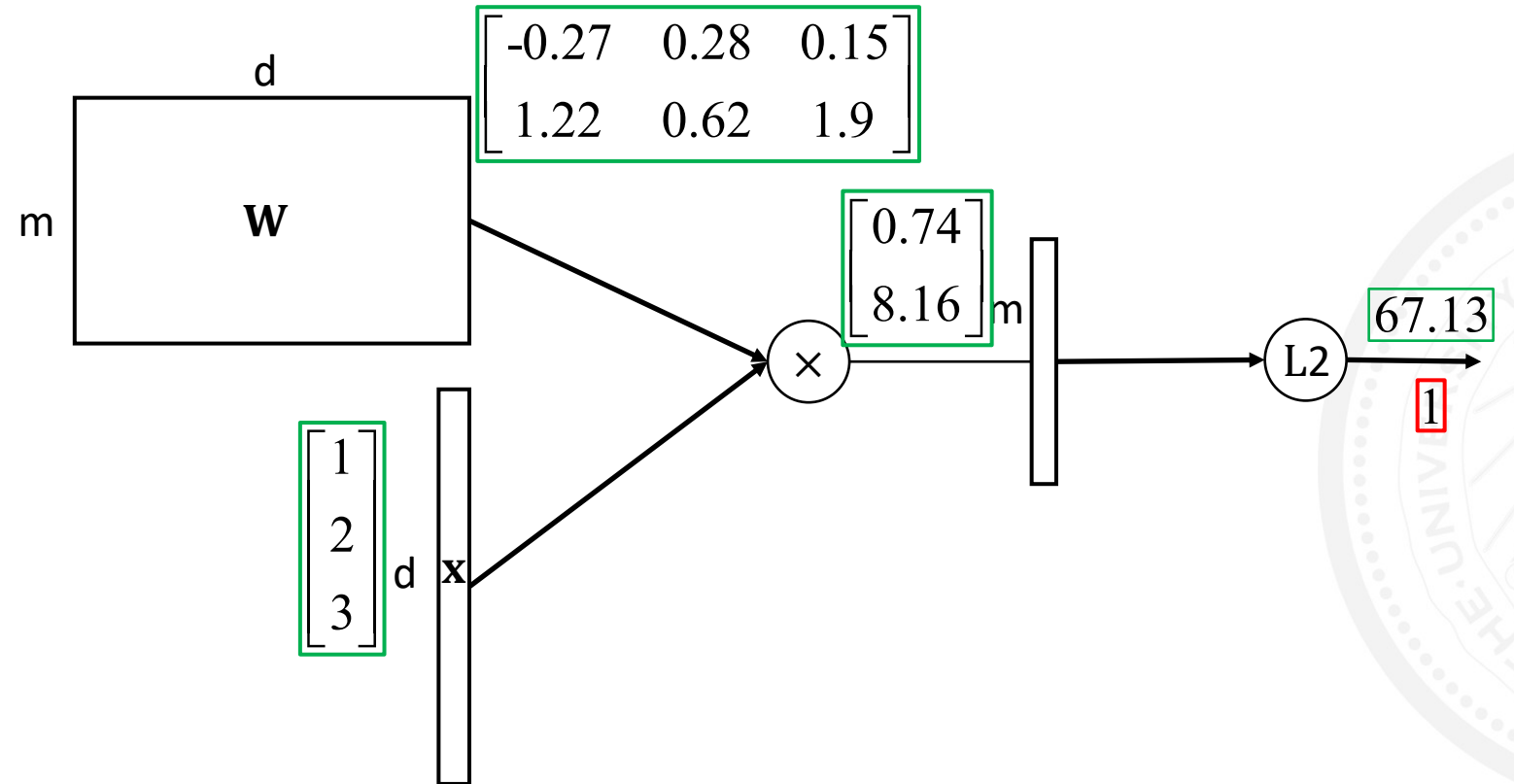




# Backpropagation

Vectorized example

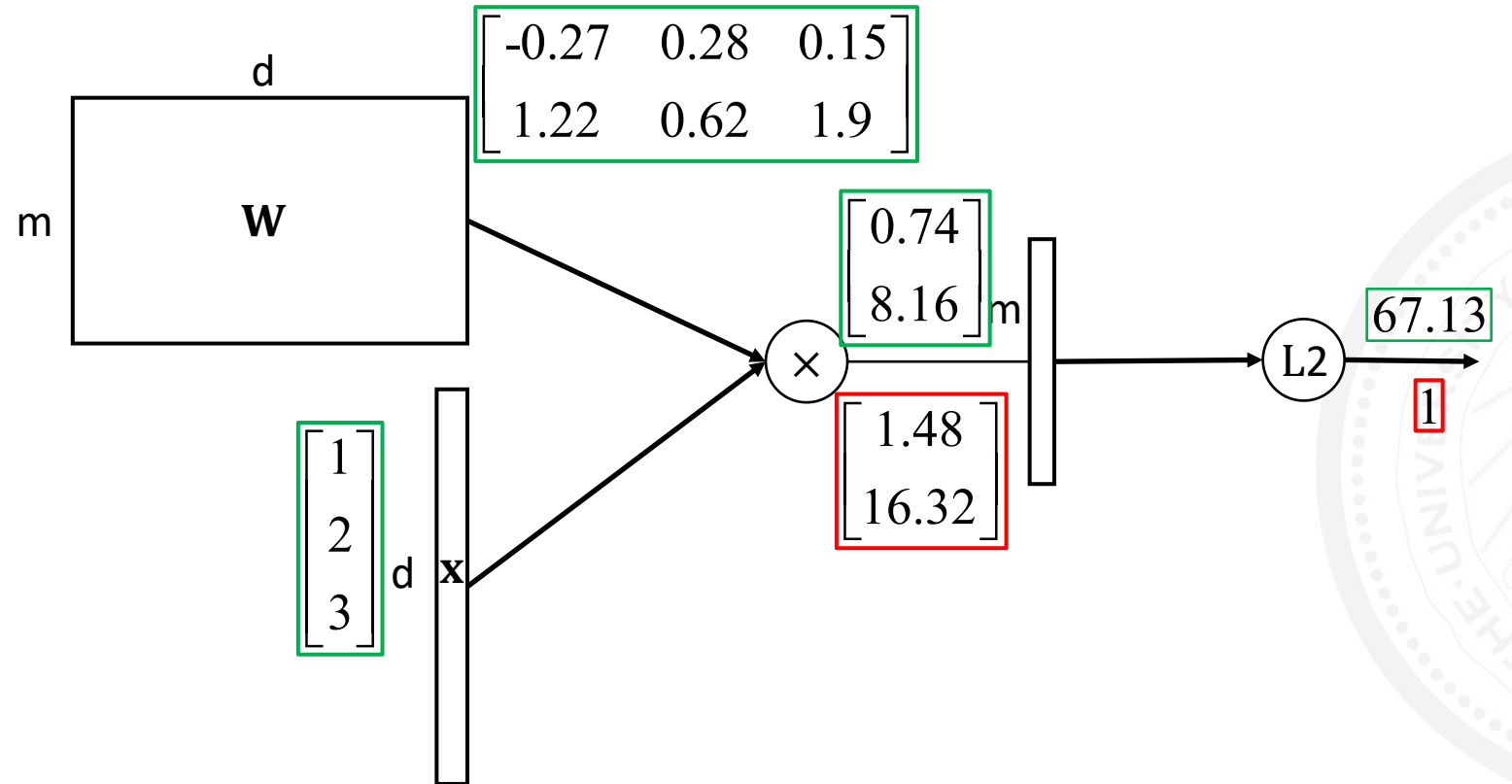
$$f(\mathbf{x}, \mathbf{W}) = \|\mathbf{W} \cdot \mathbf{x}\|^2$$



# Backpropagation

Vectorized example

$$f(\mathbf{x}, \mathbf{W}) = \|\mathbf{W} \cdot \mathbf{x}\|^2$$



# Backpropagation

Vectorized example

$$f(\mathbf{x}, \mathbf{W}) = \|\mathbf{W} \cdot \mathbf{x}\|^2$$

