



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

Learning From Data

Lecture ⁷~~6~~: Backpropagation and Neural Networks

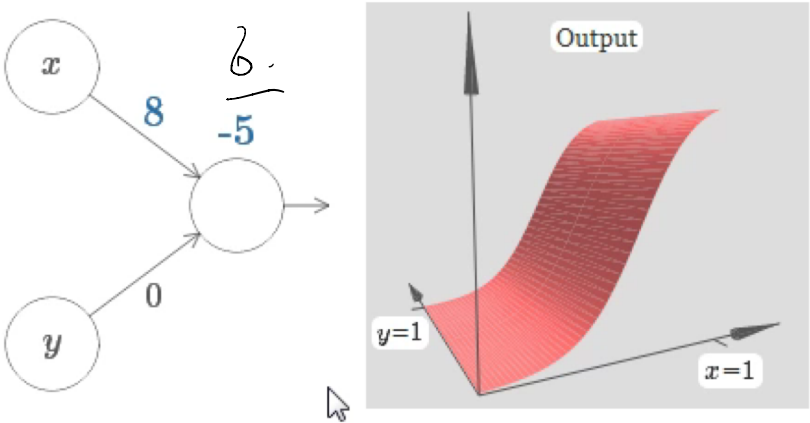
Yang Li

Slides by Lichen Wang

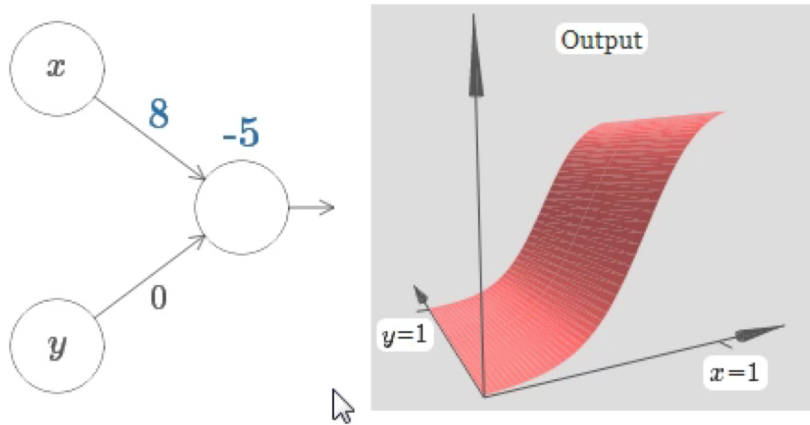
10/29/2021



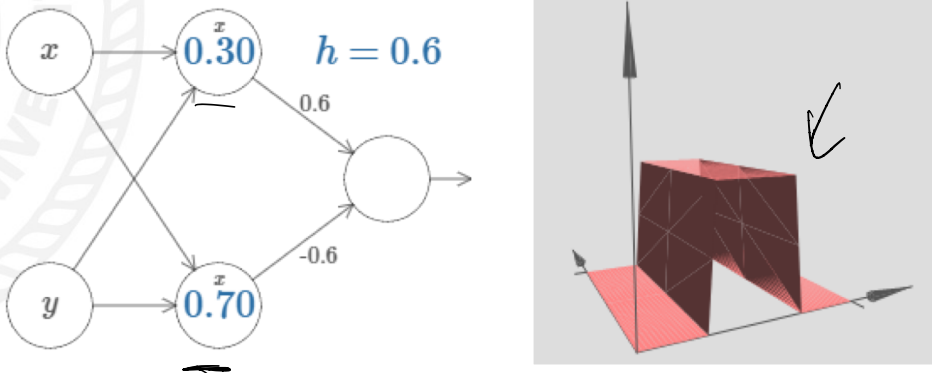
Power of single neuron



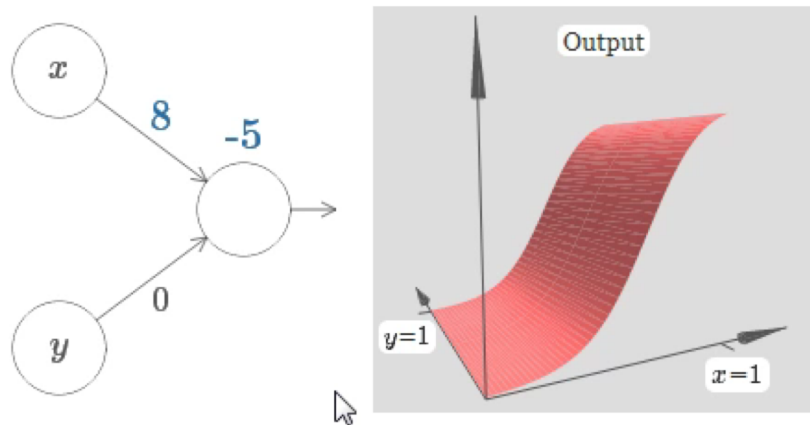
Power of single neuron



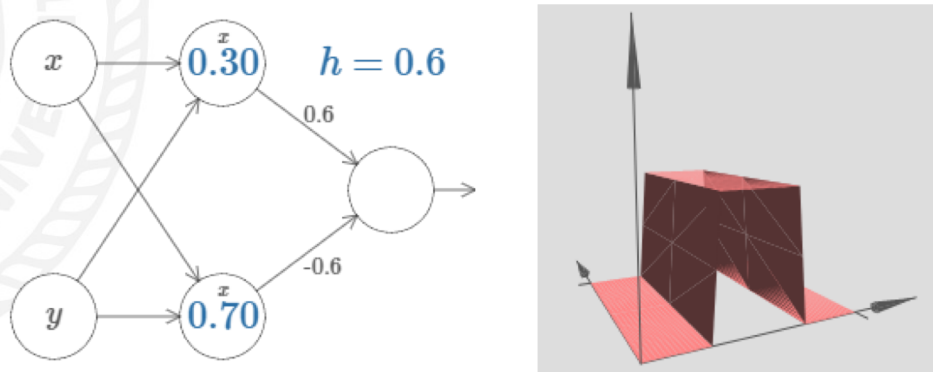
Two hidden units



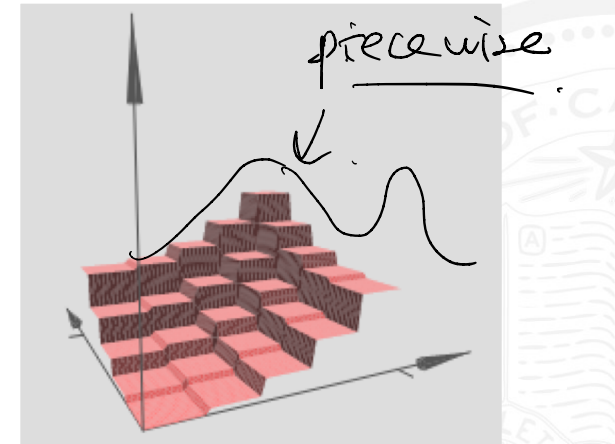
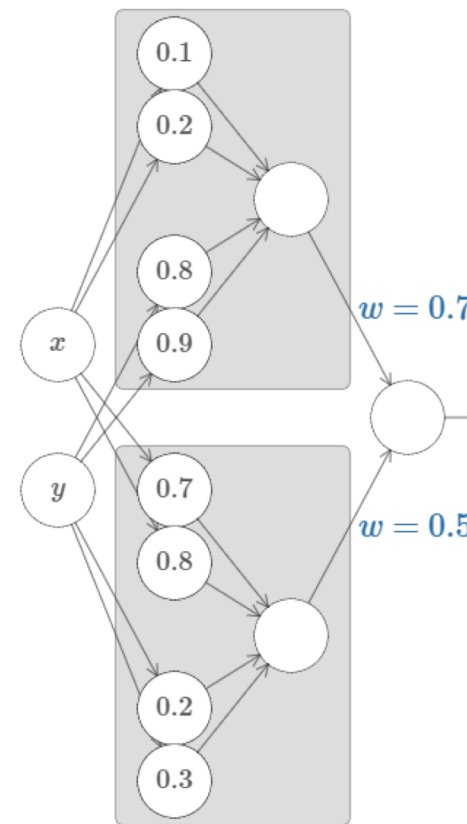
Power of single neural



Two hidden units

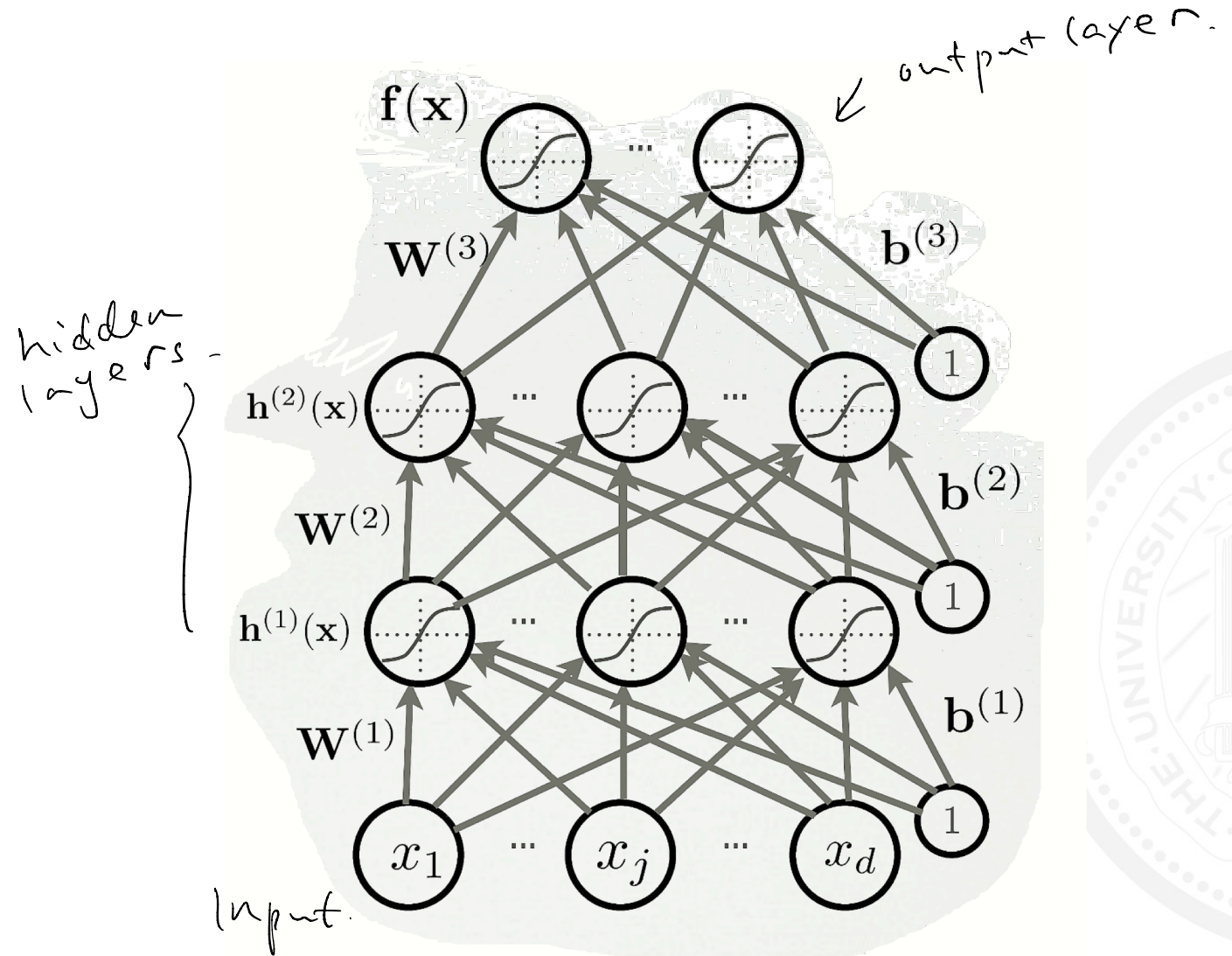


Many hidden units



Multilayer Neural Network

Could have L hidden layers

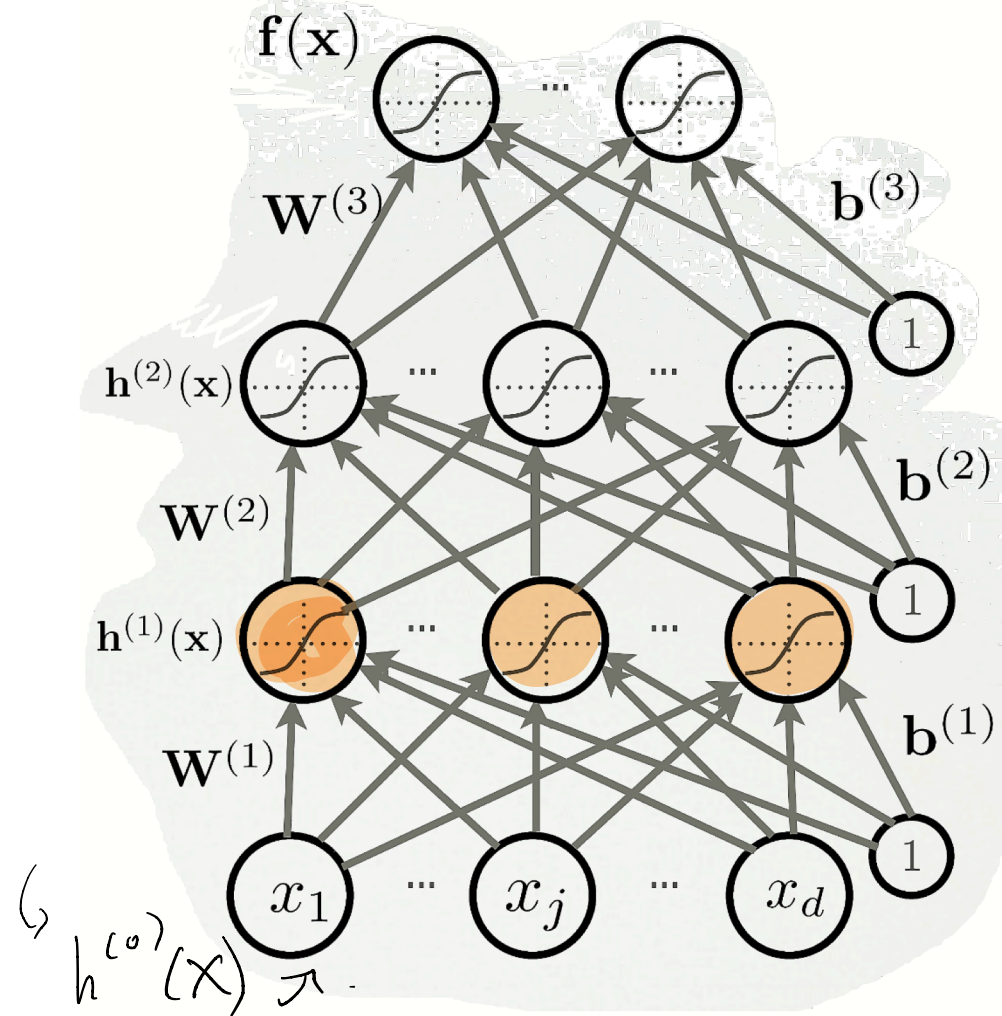


Multilayer Neural Network

Could have L hidden layers
← k th layer.

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

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



Multilayer Neural Network

Could have L hidden layers

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

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

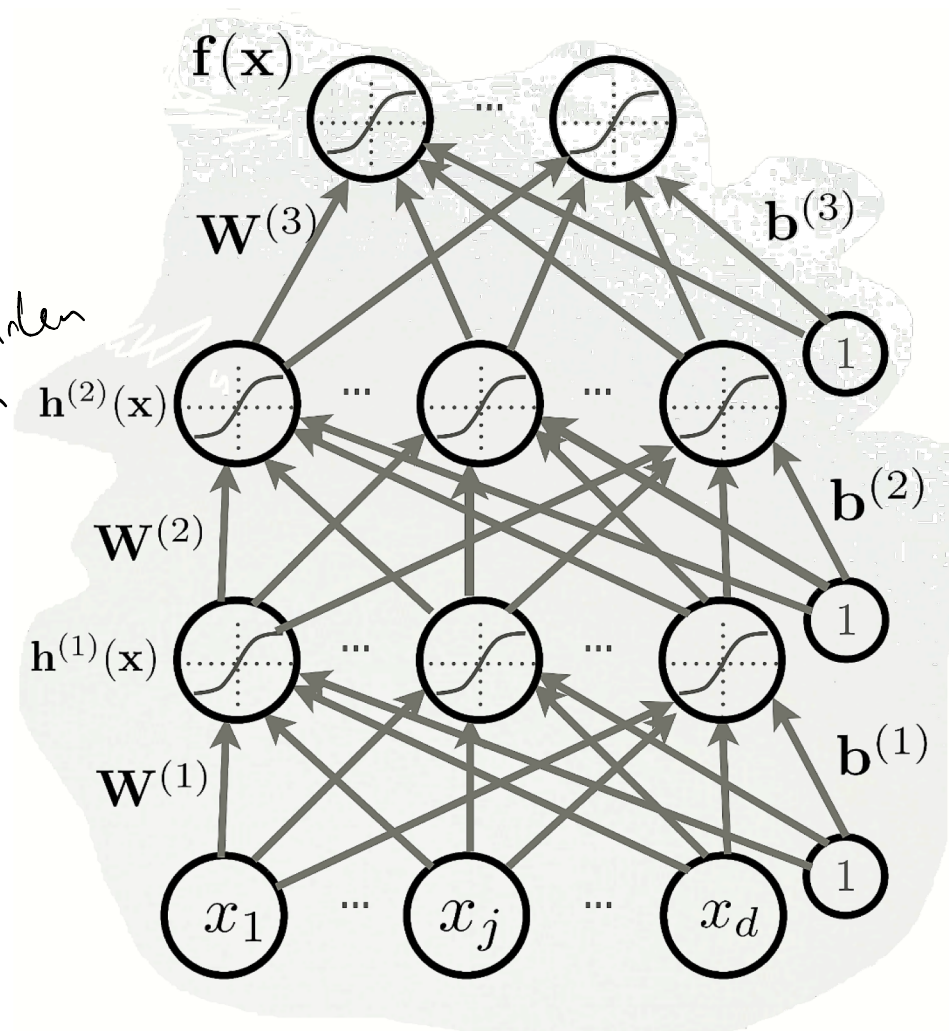
- hidden layer activation for $1 \leq k \leq L$

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

activation

Repeat
L times

one hidden
layer



Multilayer Neural Network



Could have L hidden layers

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

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

linear

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

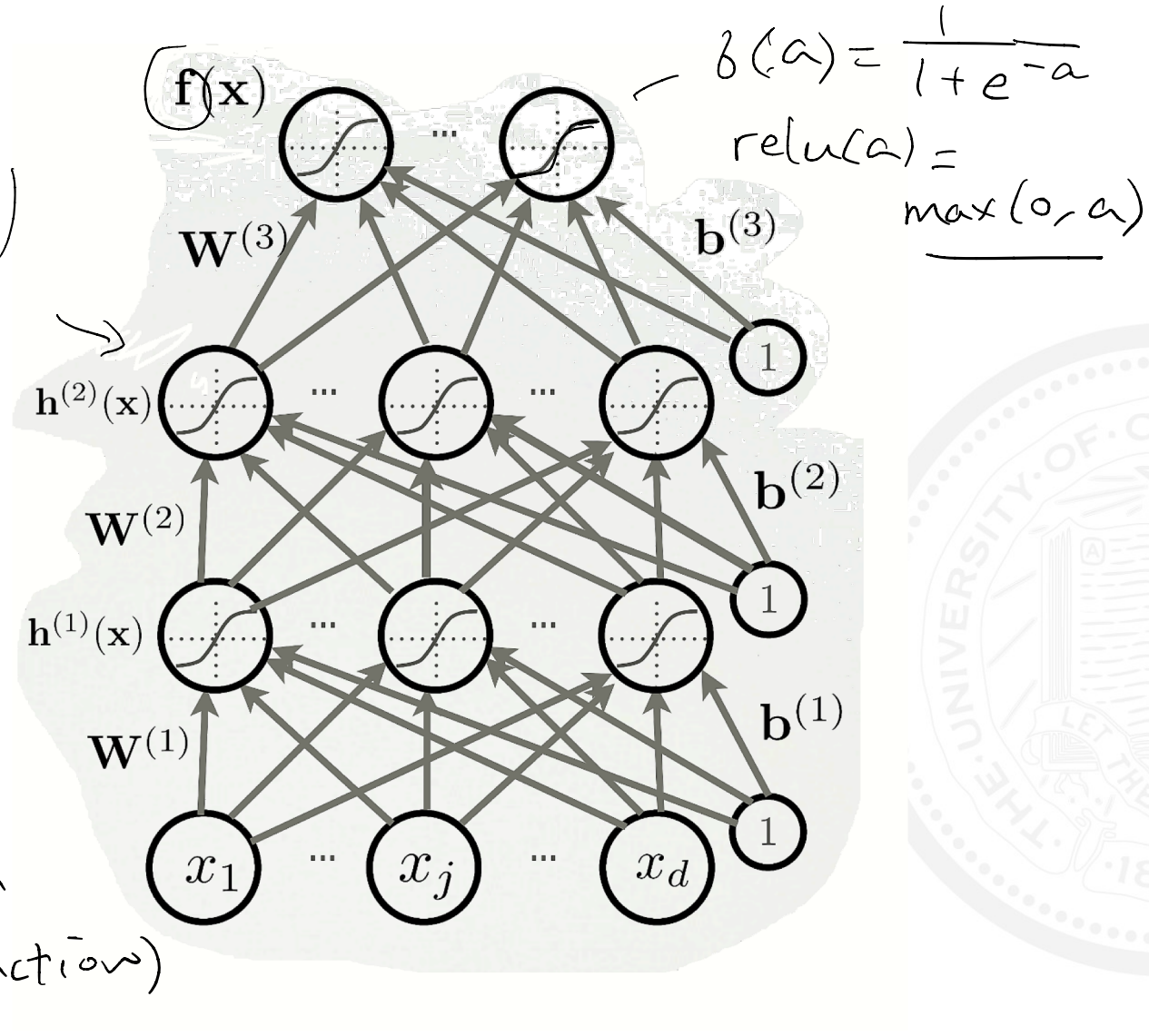
$$\mathbf{h}^{(k)}(\mathbf{x}) = \mathbf{g}(\mathbf{a}^{(k)}(\mathbf{x}))$$

threshold function (non linear)

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

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

Output function (prediction functions)



Empirical risk

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

output. ith label.

- $L(f(x^{(i)}; W), y^{(i)})$ is the loss function
- $\lambda \Omega(W)$ is the regularizer

Empirical risk

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

- $L(f(x^{(i)}; W), y^{(i)})$ is the loss function
- $\lambda \Omega(W)$ is the regularizer

Softmax loss for sample $(x^{(i)}, y^{(i)})$

$$L_i = -\log \left(\frac{e^{f_{y^{(i)}}}}{\sum_{j \geq 1} e^{f_j}} \right)$$

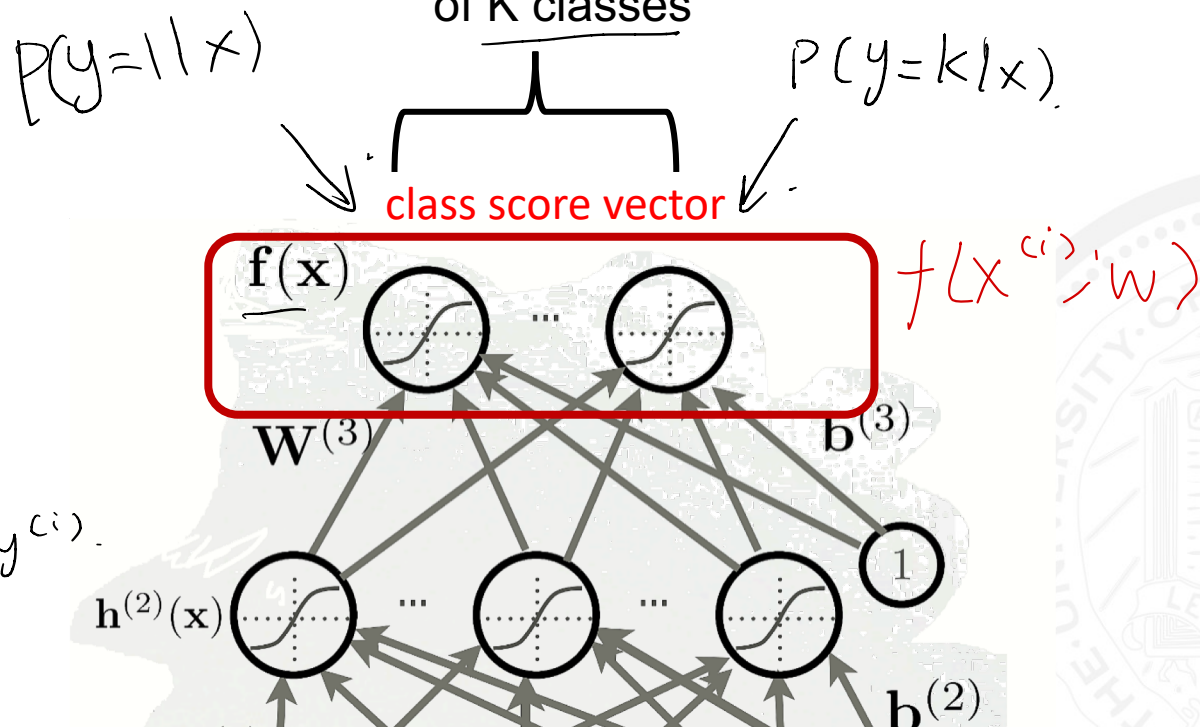
$\in \{1, \dots, K\}$



f_j : = j -th element of class score vector $f(x^{(i)}; W)$

Softmax example:

Unnormalized class probability
of K classes



Empirical risk

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



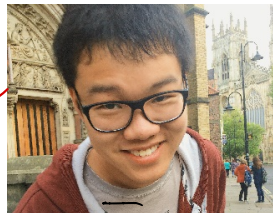
- $L(f(x^{(i)}; W), y^{(i)})$ is the loss function
- $\lambda \Omega(W)$ is the regularizer

Softmax loss for sample $(x^{(i)}, y^{(i)})$

$$L_i = -\log \left(\frac{e^{f_{y^{(i)}}}}{\sum_j e^{f_j}} \right)$$

f_j : = j -th element of class score vector $f(x^{(i)}; W)$

Softmax example: f_1 f_2 f_3 Loss

$x^{(i)}$	Beckham	Yanzu	Lichen	L_i $f_{y^{(i)}}$
	4.9	1.1	-0.9	$-\log \left(\frac{e^{4.9}}{e^{4.9} + e^{1.1} + e^{-0.9}} \right)$
	-2.6	1.7	1.2	$-\log \left(\frac{e^{1.7}}{e^{-2.6} + e^{1.7} + e^{1.2}} \right)$
	0.2	1.2	2.2	$-\log \left(\frac{e^{2.2}}{e^{0.2} + e^{1.2} + e^{2.2}} \right)$

$f(x^{(i)}; W)$
given some weight \overline{W} .

how to evaluate loss
given (x_i, y_i, W)



- Find the optimal parameter

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

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

Find the optimal parameter

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

Stochastic Gradient Descent (SGD)

Algorithm

1. Initialize W

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

2a. $\Delta := -(\nabla_w L(f(x^{(i)}; W), y^{(i)}) + \lambda \nabla_w \Omega(w))$ *gradient update (negative gradient direction)*

2b. $W \leftarrow W + \alpha \Delta$ *learning rate $\alpha > 0$.*

Training epoch = Iterating over all examples

Handwritten notes: $\lambda \geq 0$.

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

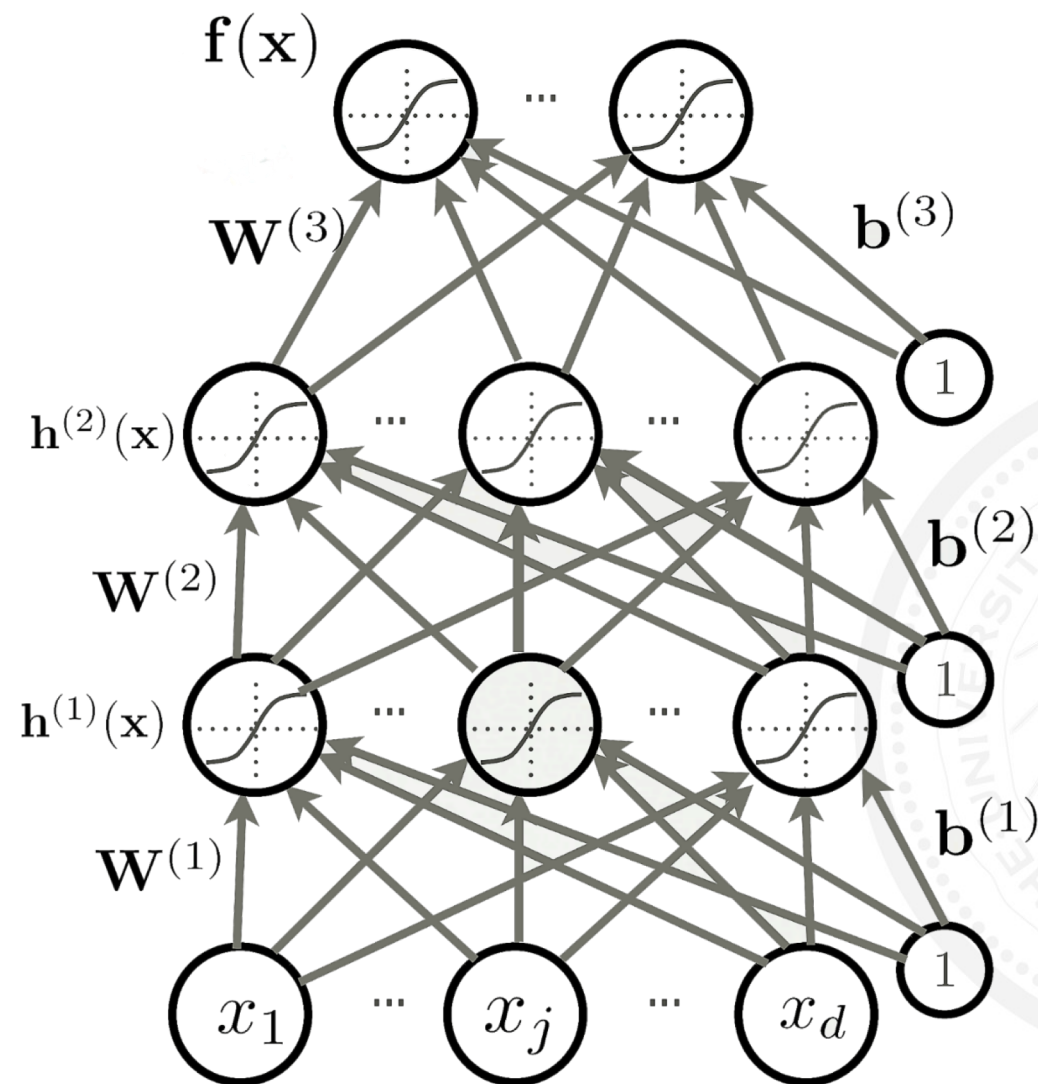
Animation: (for SVM loss)

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

Follow the slope



How many parameter do we have?

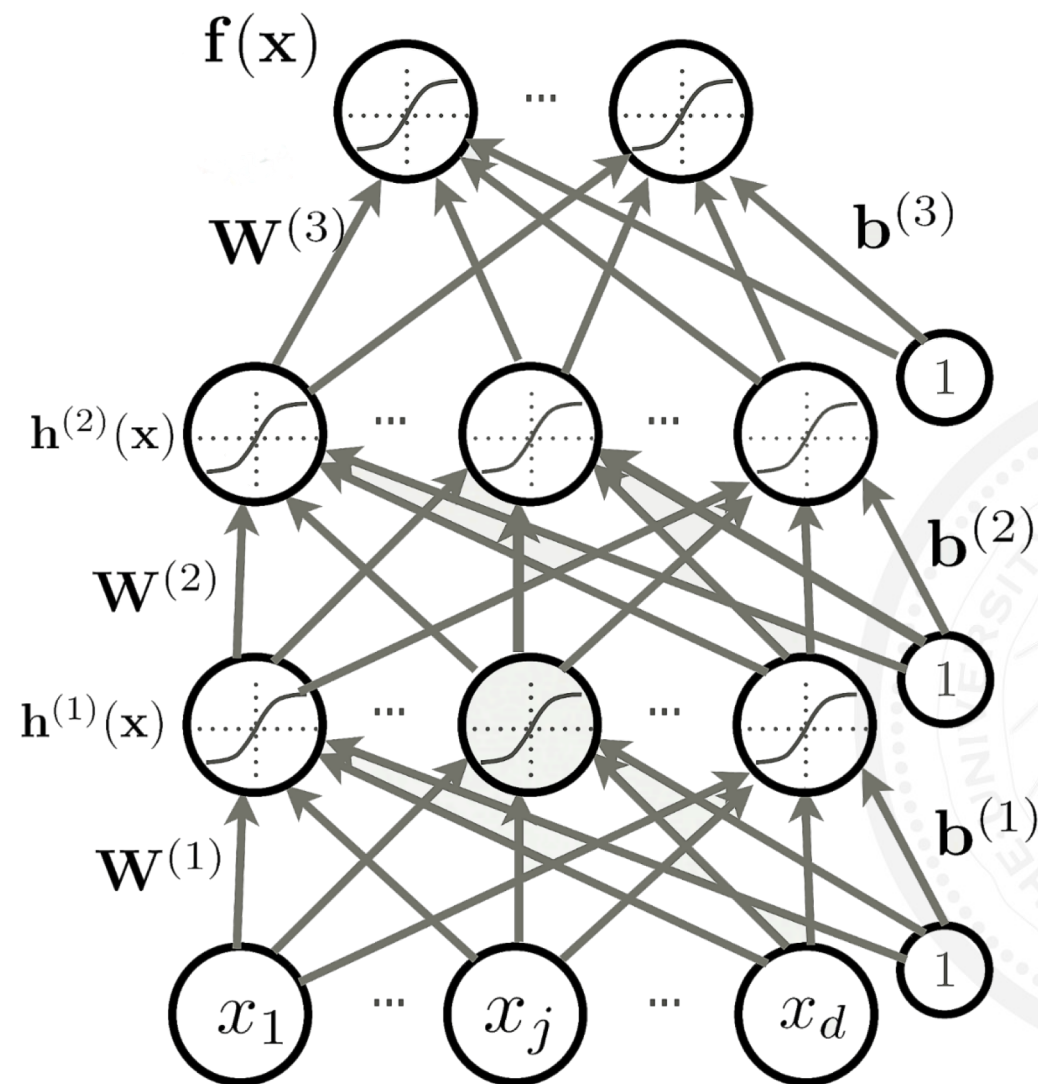


Follow the slope



How many parameter do we have?

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



Follow the slope



TBSI

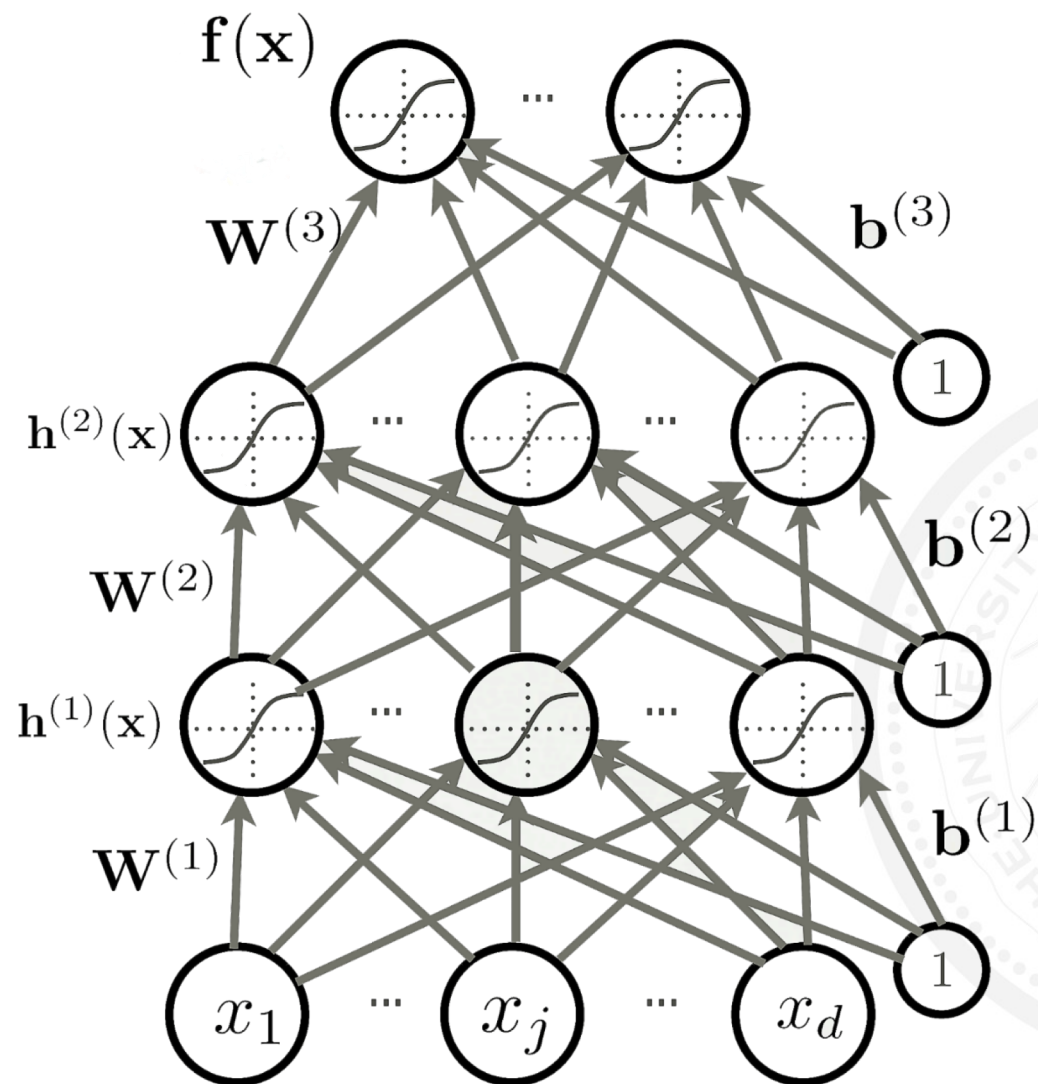
清华-伯克利深圳学院
Tsinghua-Berkeley Shenzhen Institute

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

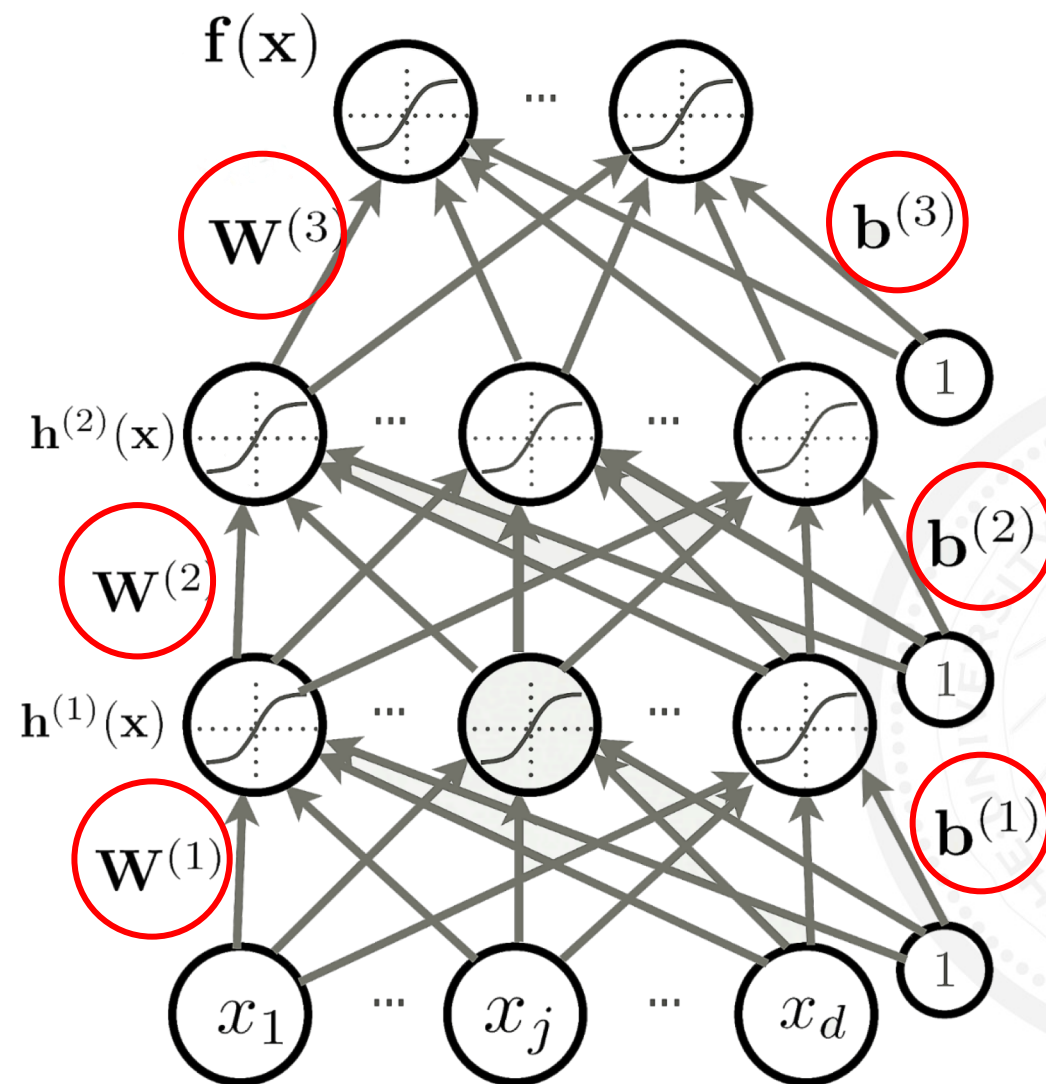


How many parameter do we have?

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

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

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):

[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

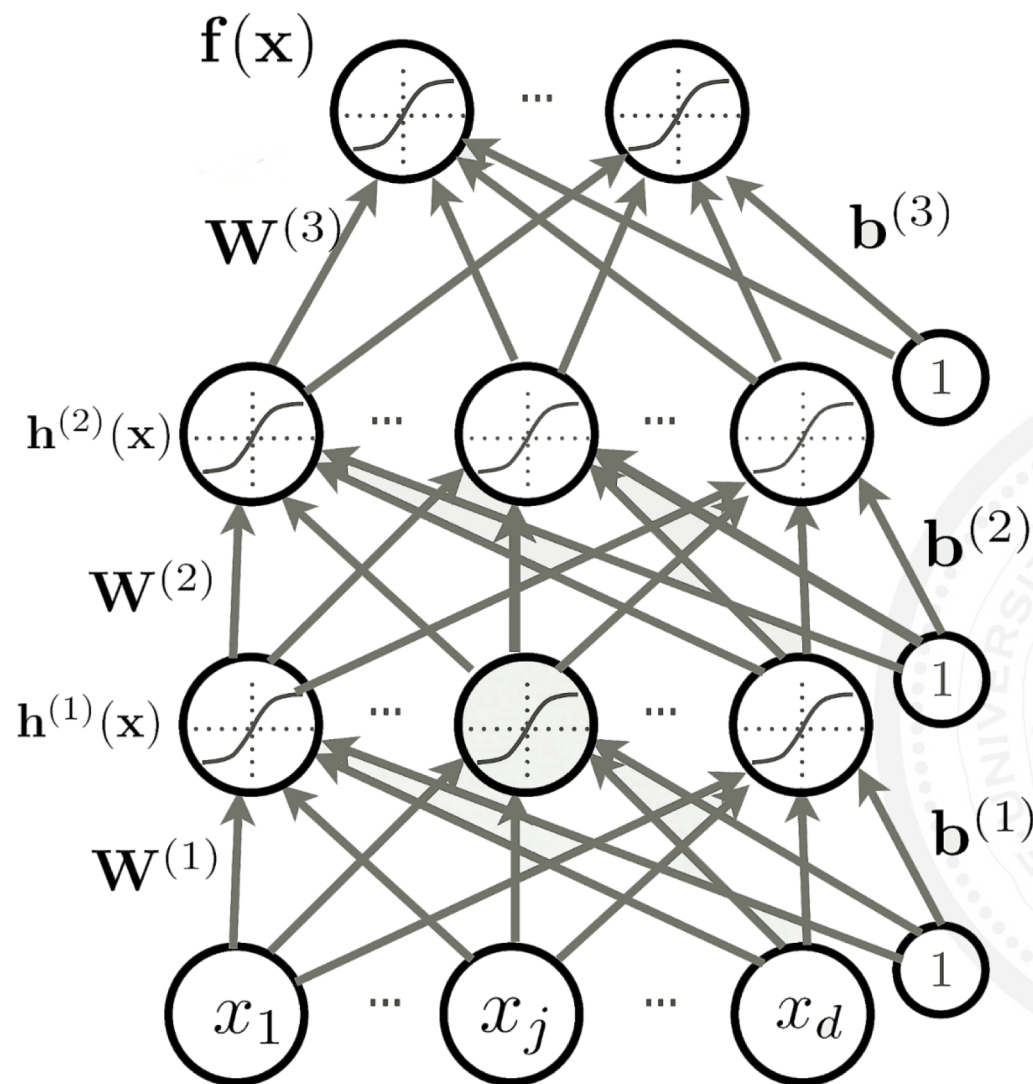


How many parameter do we have?

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



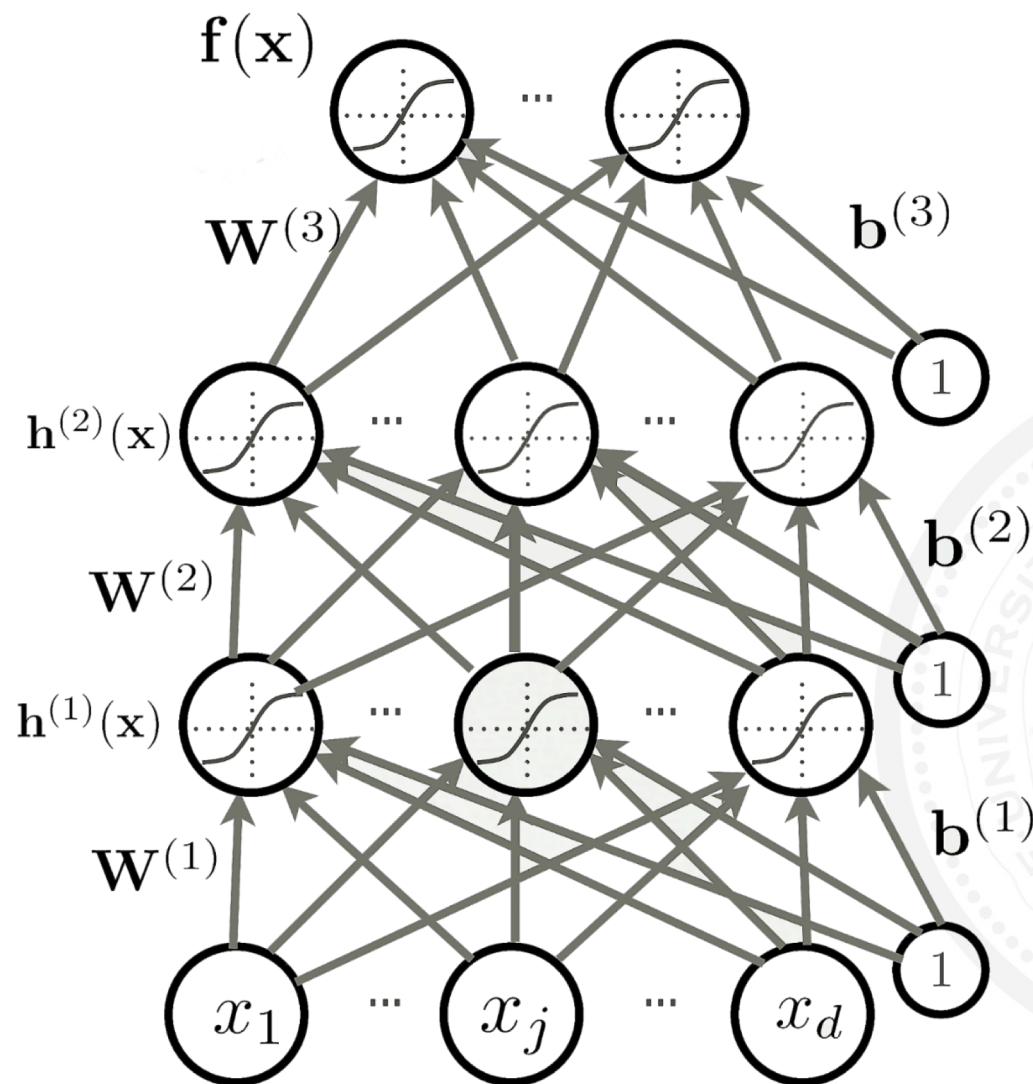
How many parameter do we have?

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!

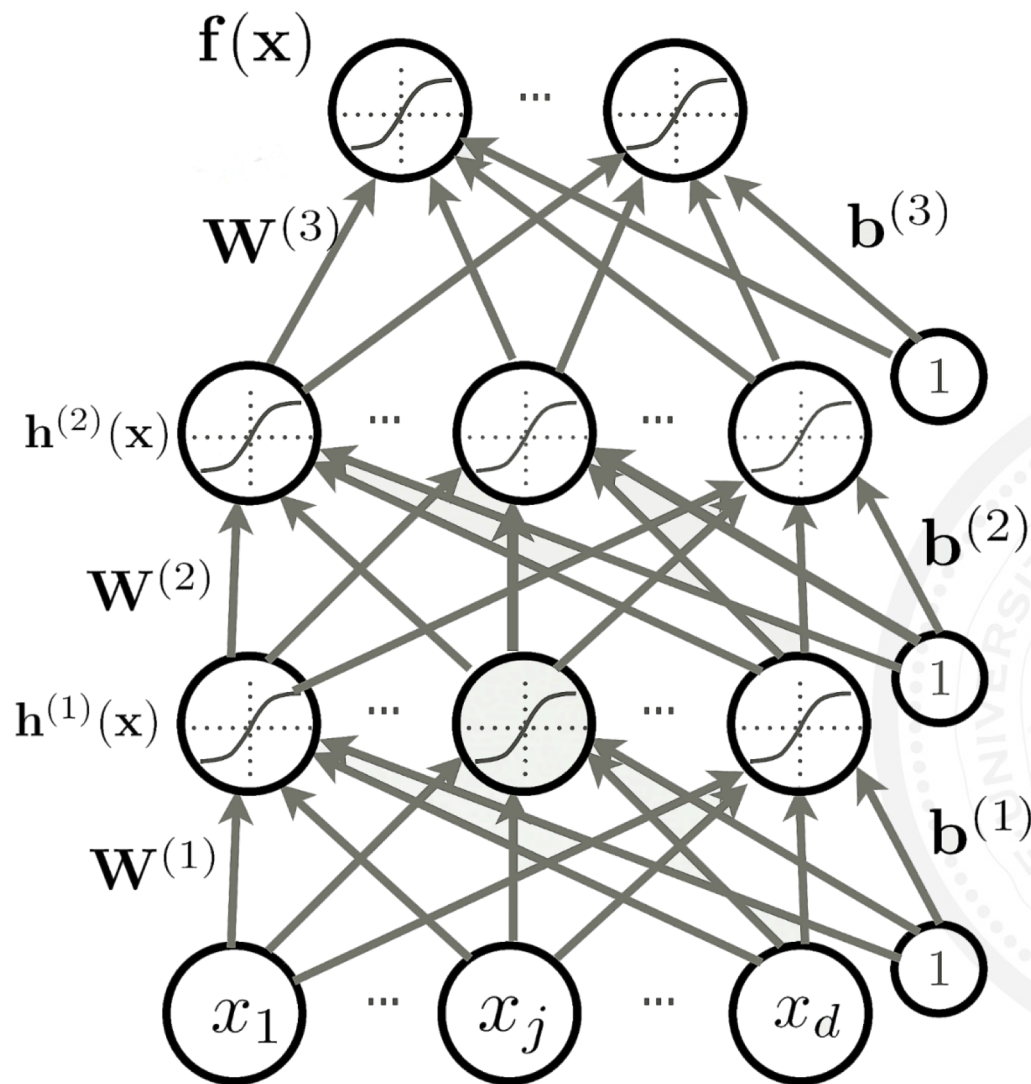
$$\hat{L} = \frac{1}{N} \sum_i (L_i(f(\mathbf{x}^{(i)}; \mathbf{W}), y^{(i)}) + \lambda \Omega(\mathbf{W}))$$

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

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

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

Analytic gradient: exact, fast, error-prone



Backpropagation



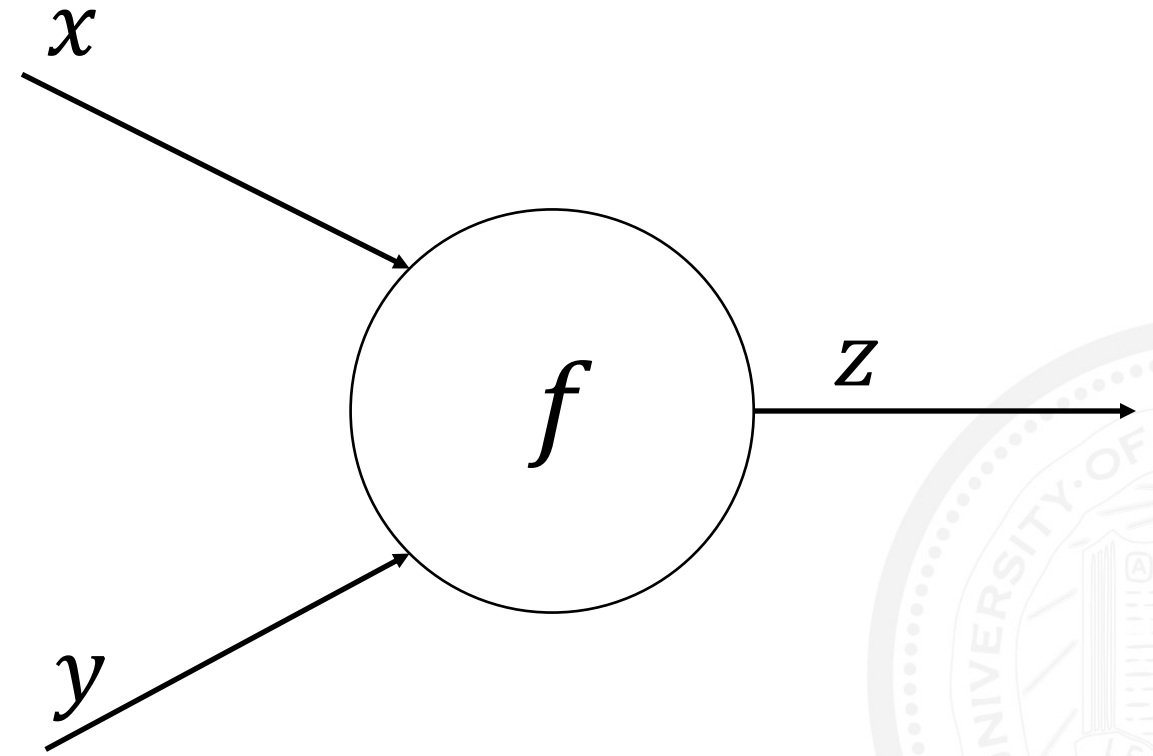
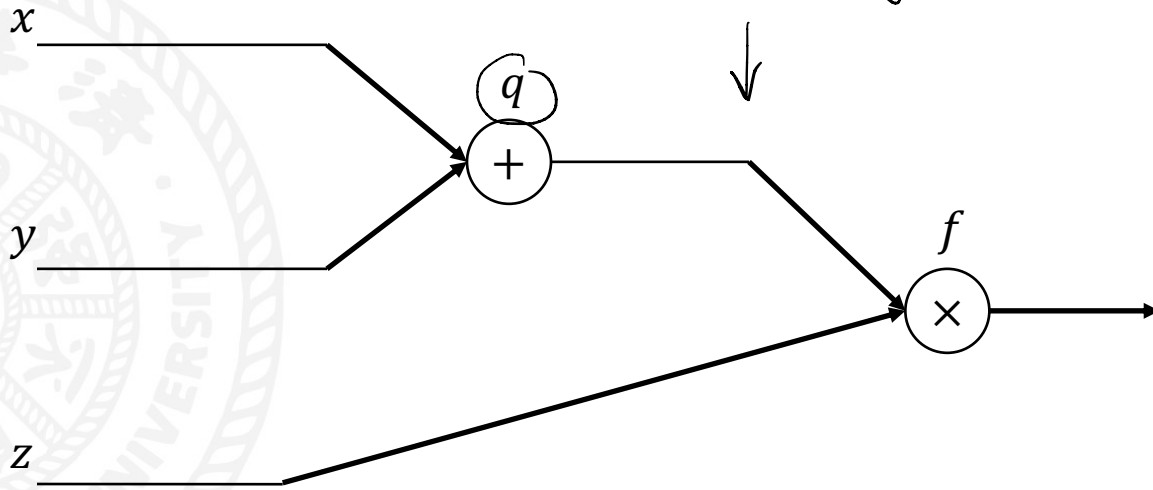
TBSI

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

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

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

computation graph.

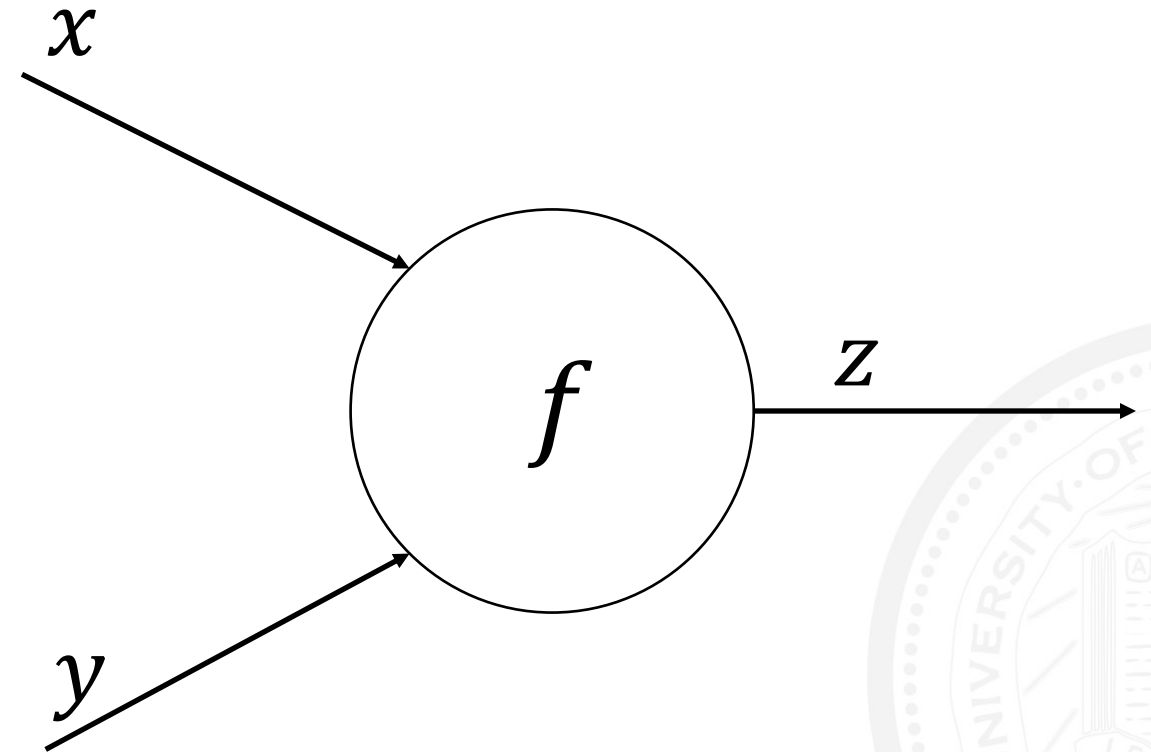
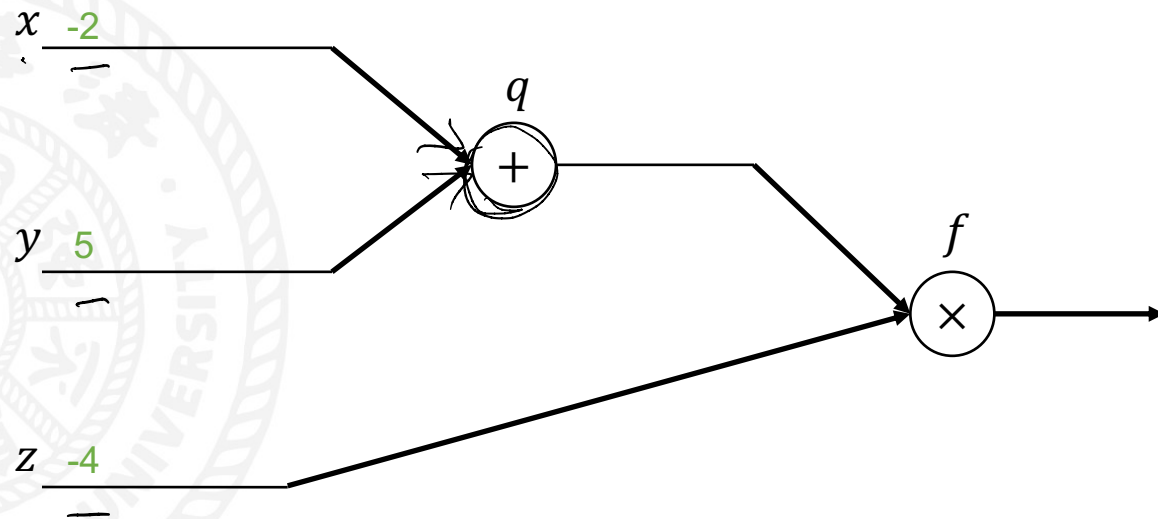


Backpropagation



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

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



Backpropagation

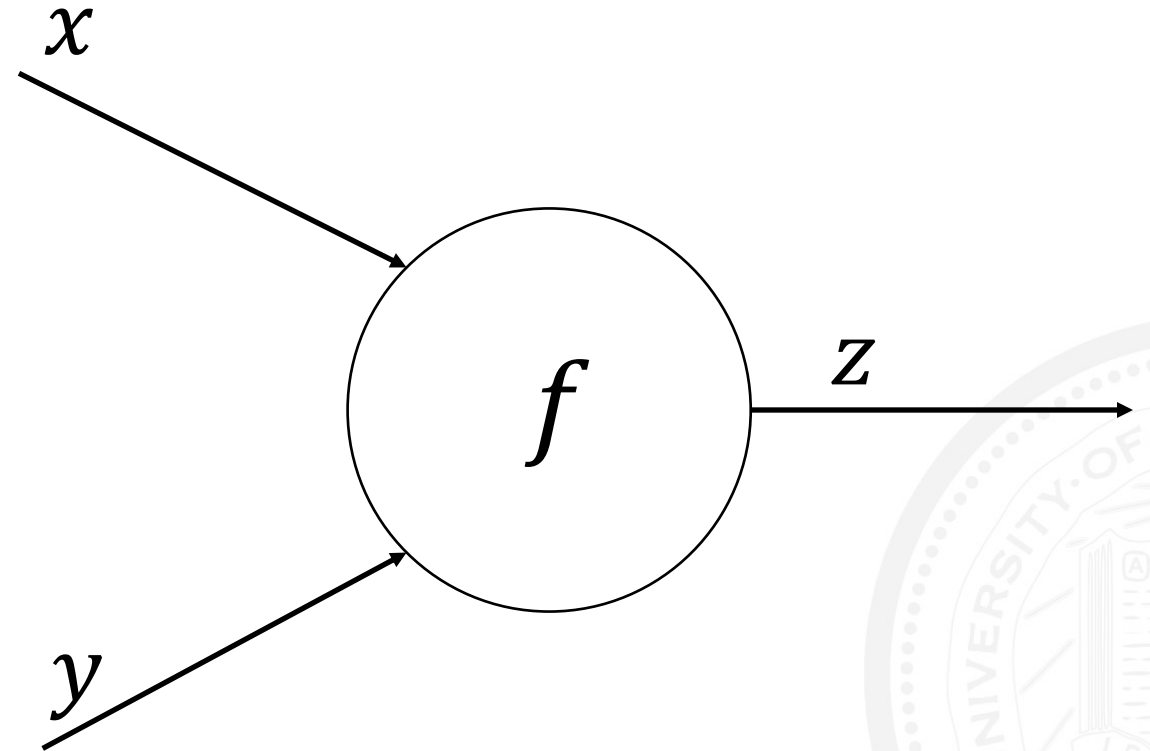
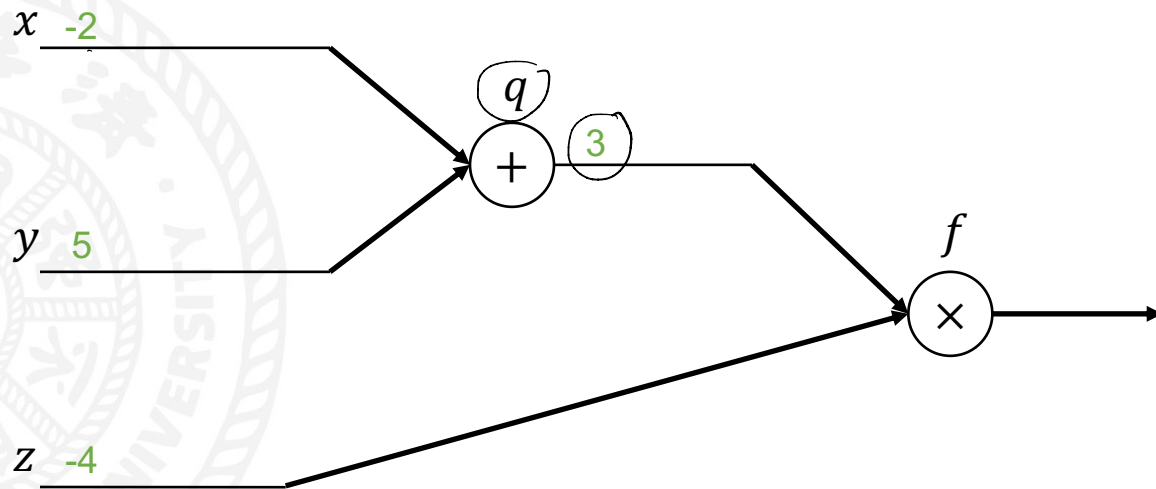


TBSI

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

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

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



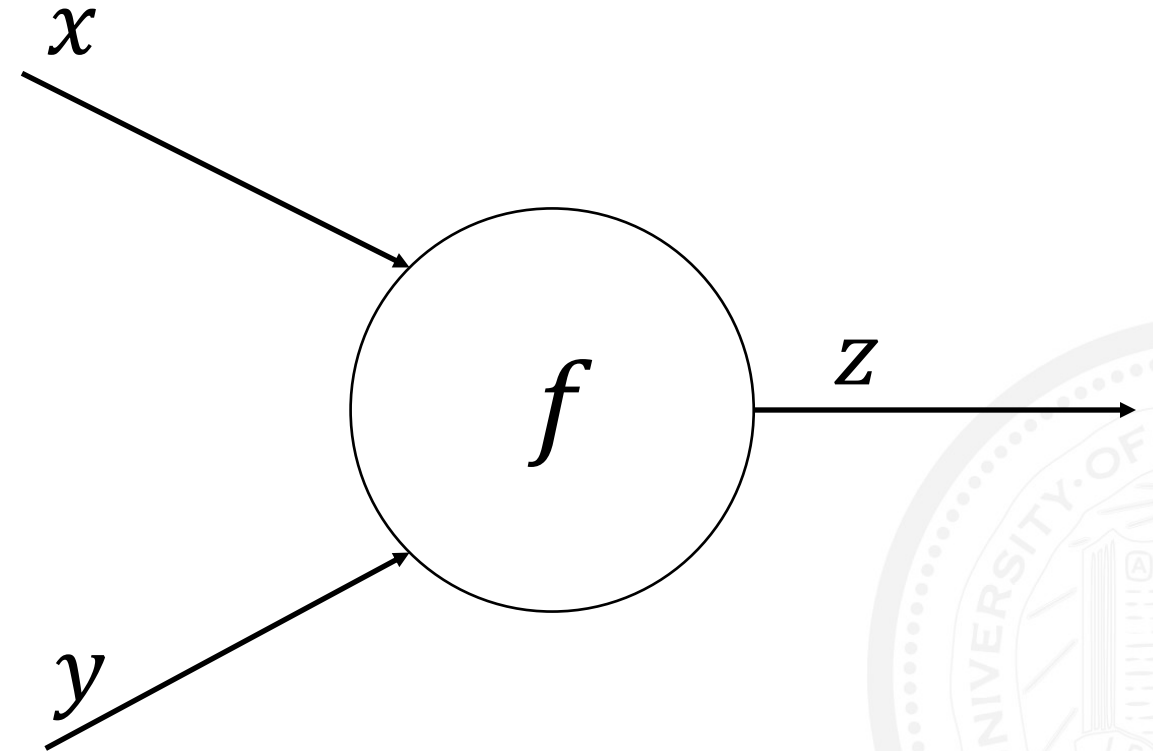
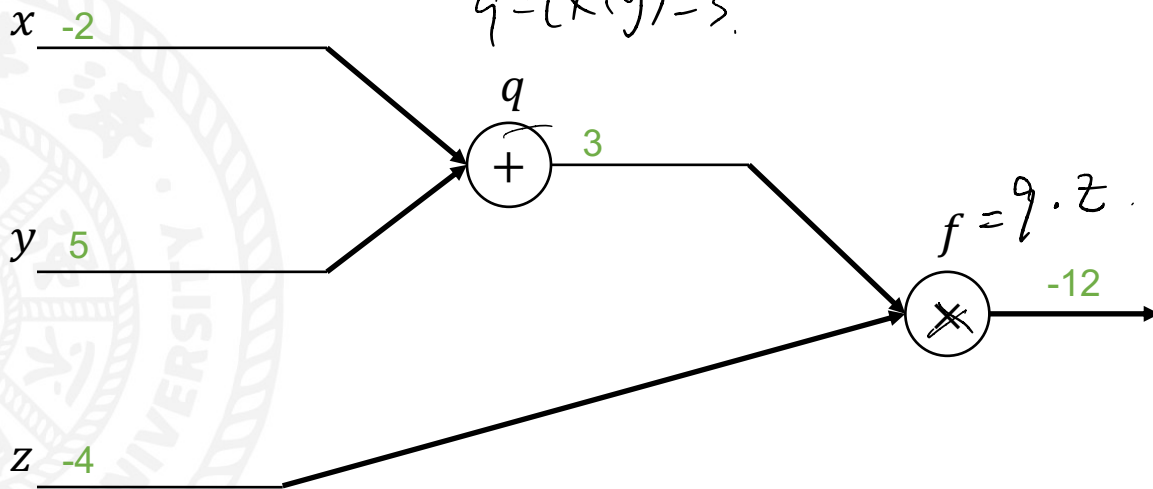
Backpropagation



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

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

$$q = (x + y) = 3$$

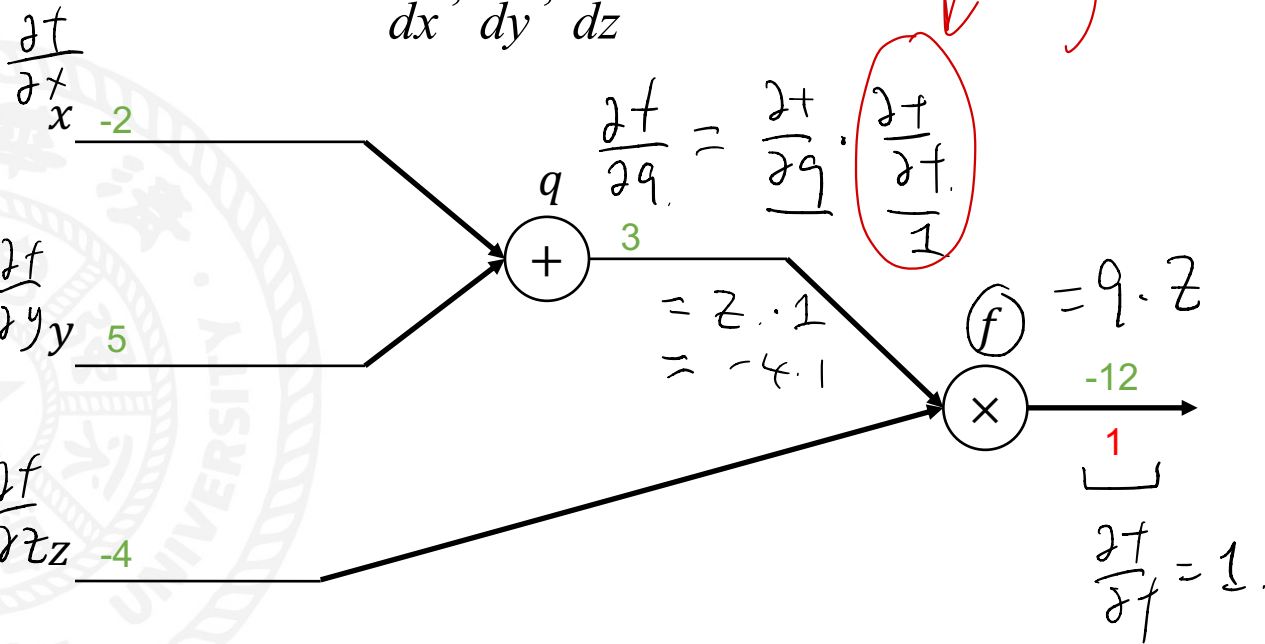


Backpropagation

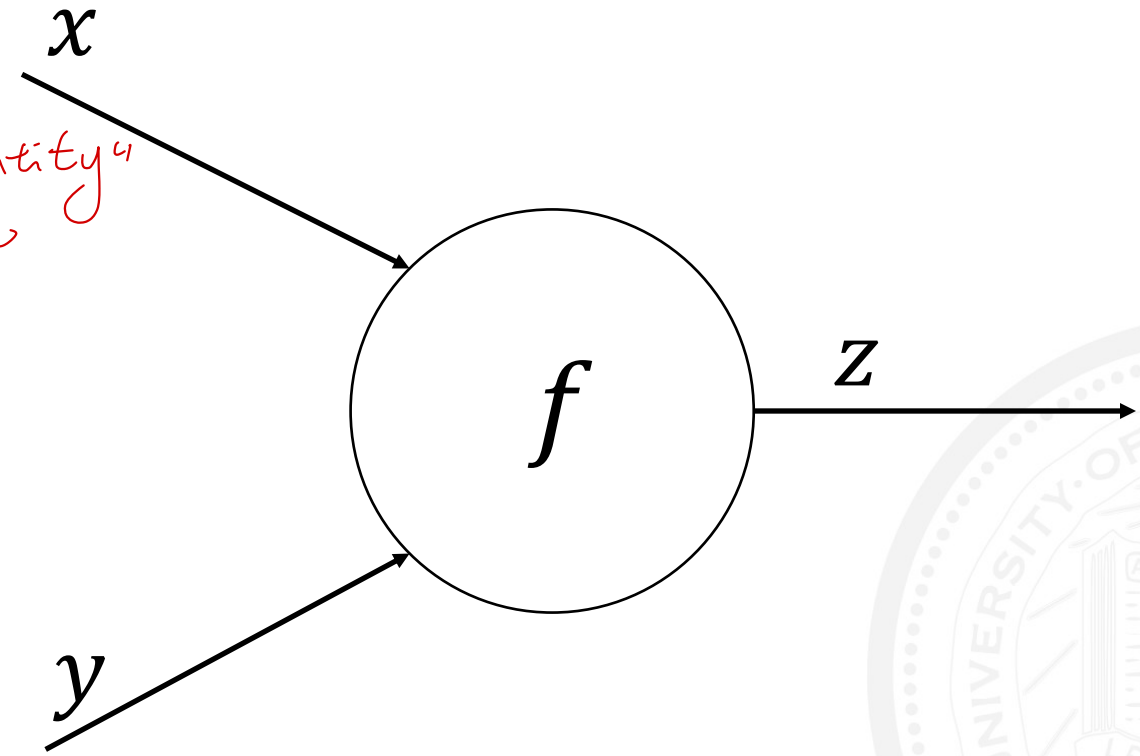


$$f(x, y, z) = (x + y)z$$

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



derivative of the "identity" function

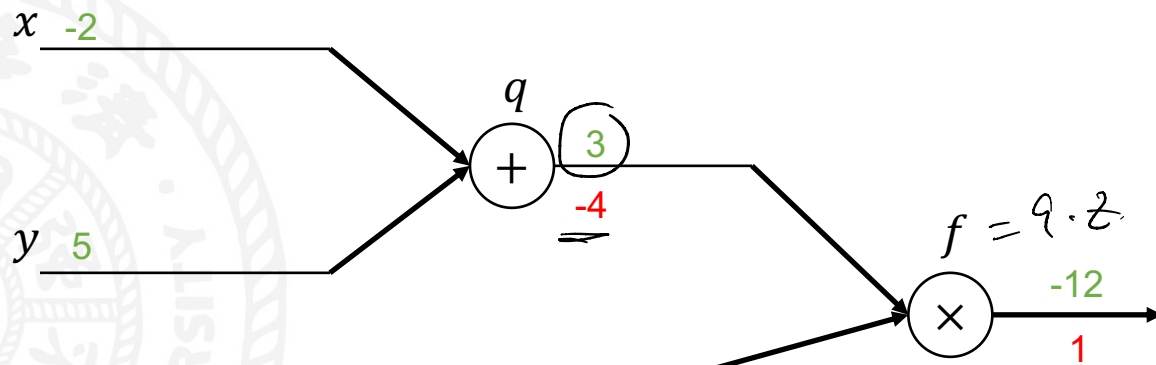


Backpropagation

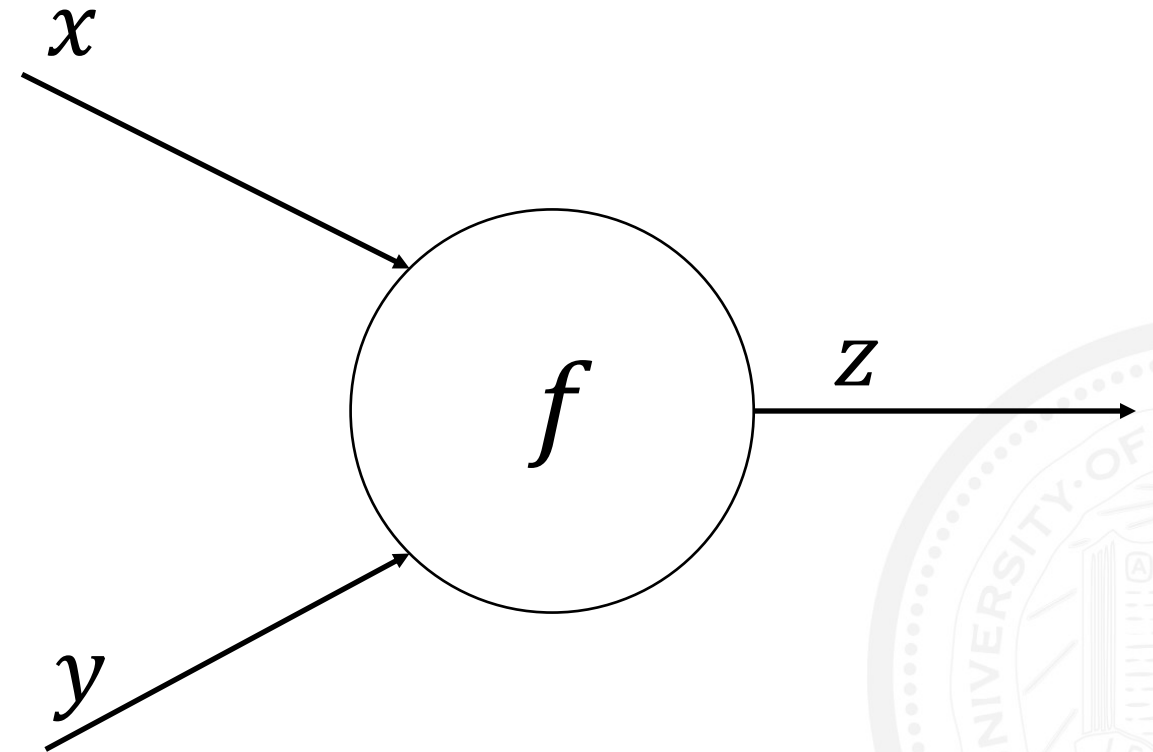


$$f(x, y, z) = (x + y)z$$

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



$$\frac{\partial f}{\partial z} = q$$
$$\frac{\partial f}{\partial z} \cdot \frac{\partial q}{\partial x} = q \cdot 1 = 3 \cdot 1$$



Backpropagation



TBSI

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

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

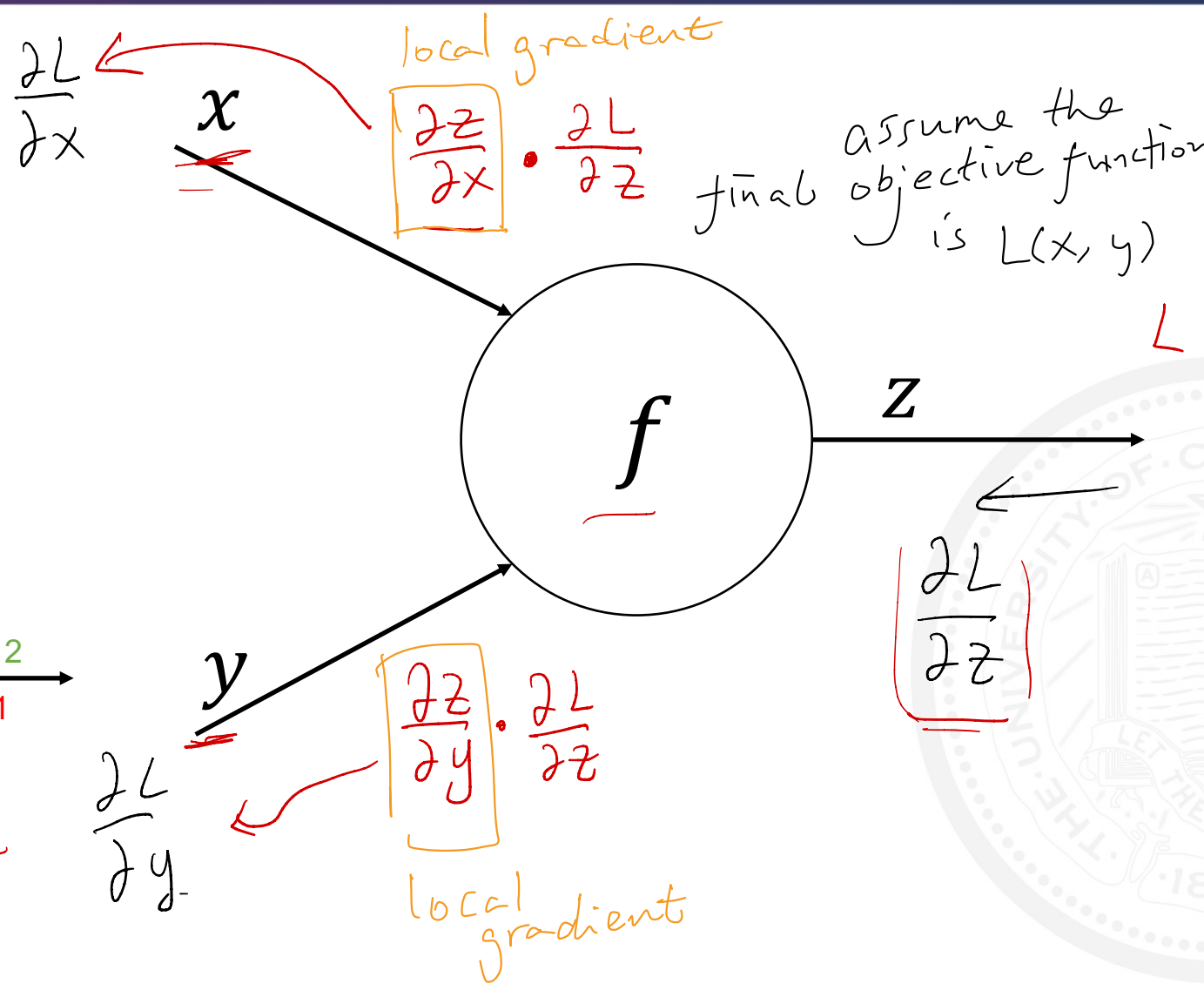
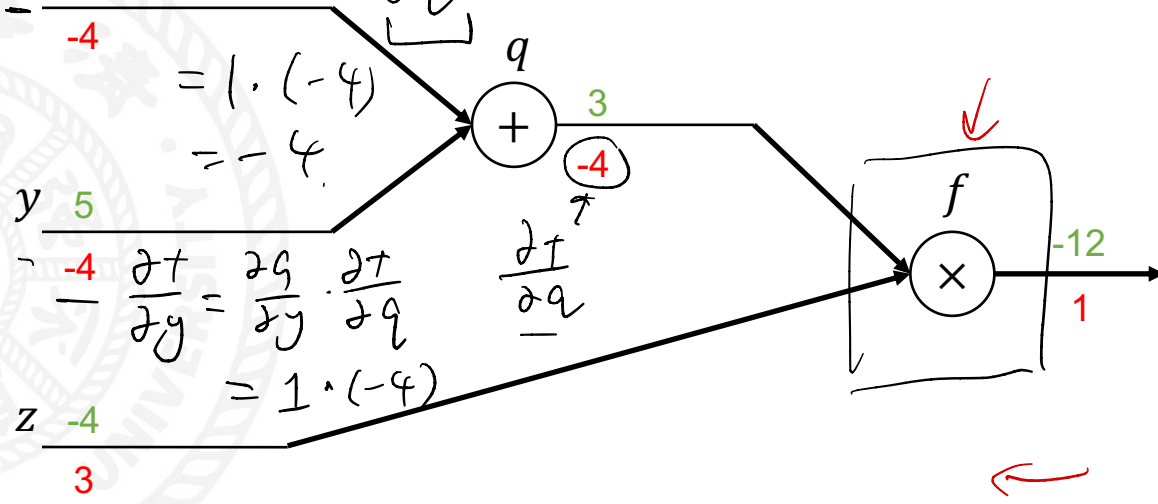
We want $\frac{df}{dx}, \frac{df}{dy}, \frac{df}{dz}$

$$\frac{\partial f}{\partial x} = \frac{\partial q}{\partial x} \cdot \frac{\partial f}{\partial q} \quad q = x + y$$

$$\frac{\partial f}{\partial x} = 1 \cdot (-4) = -4$$

$$\frac{\partial f}{\partial y} = 1 \cdot (-4) = -4$$

$$\frac{\partial f}{\partial z} = 1 \cdot (-4) = -4$$

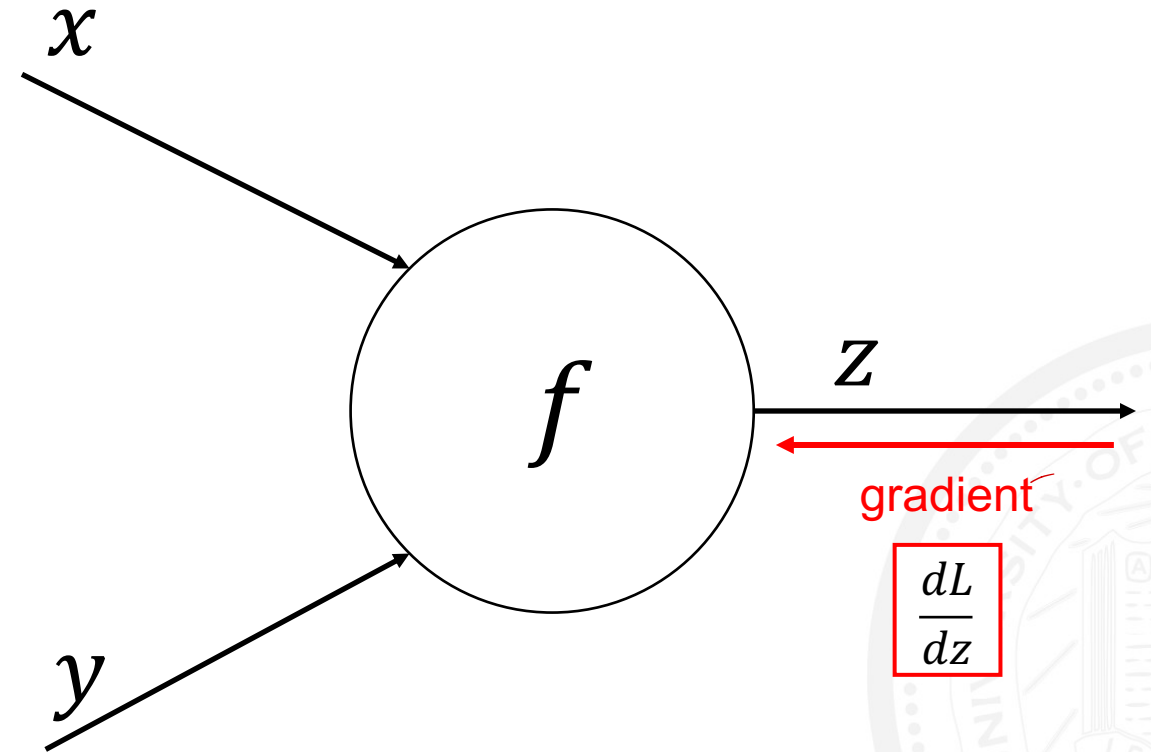
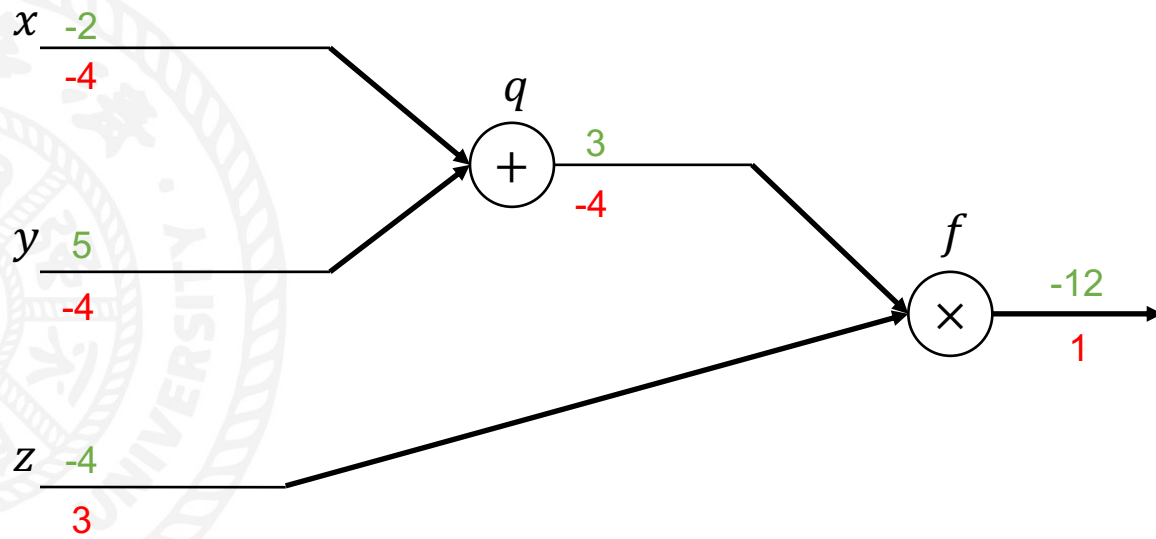


Backpropagation



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

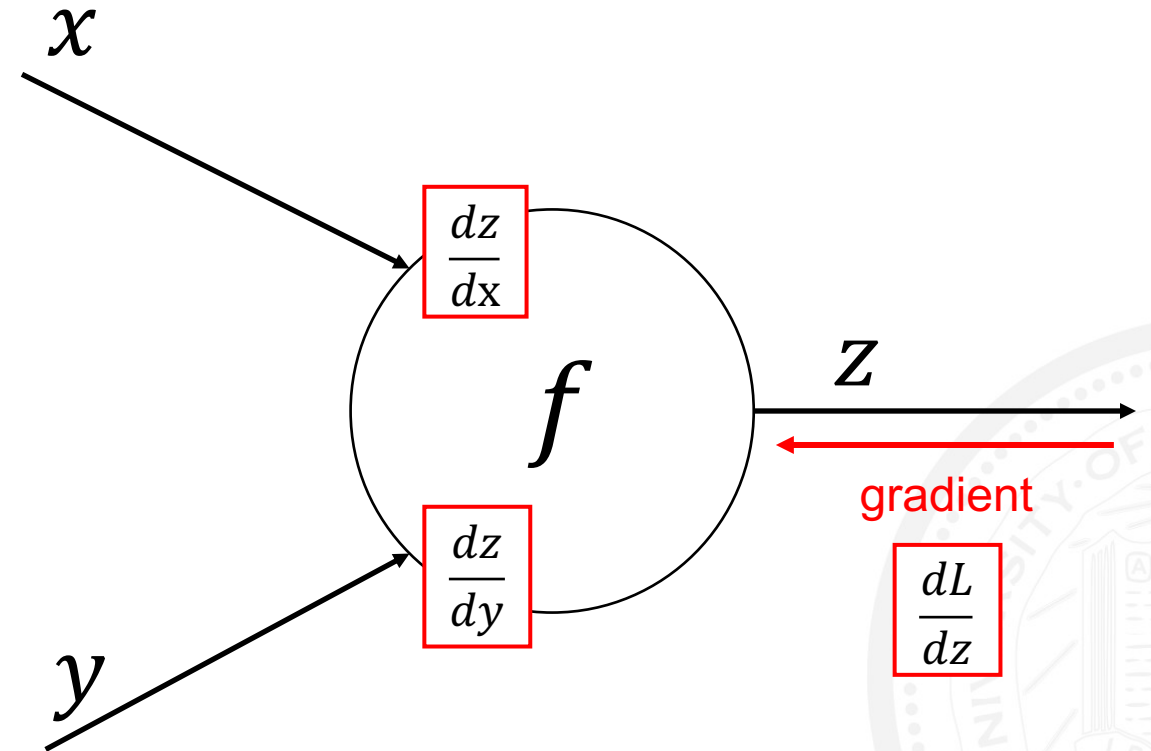
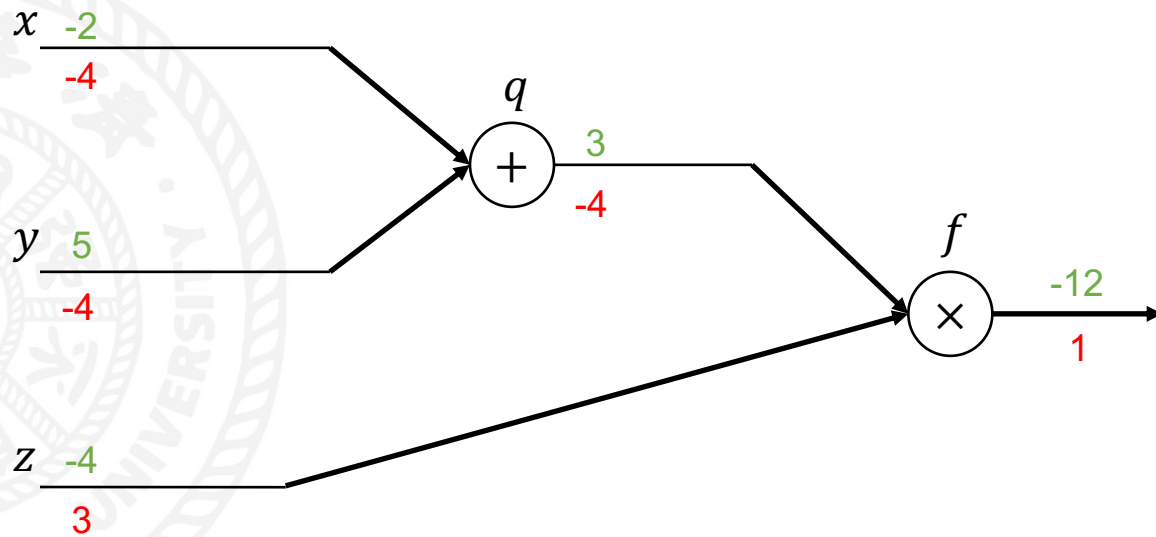
We want $\frac{df}{dx}$, $\frac{df}{dy}$, $\frac{df}{dz}$



Backpropagation

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

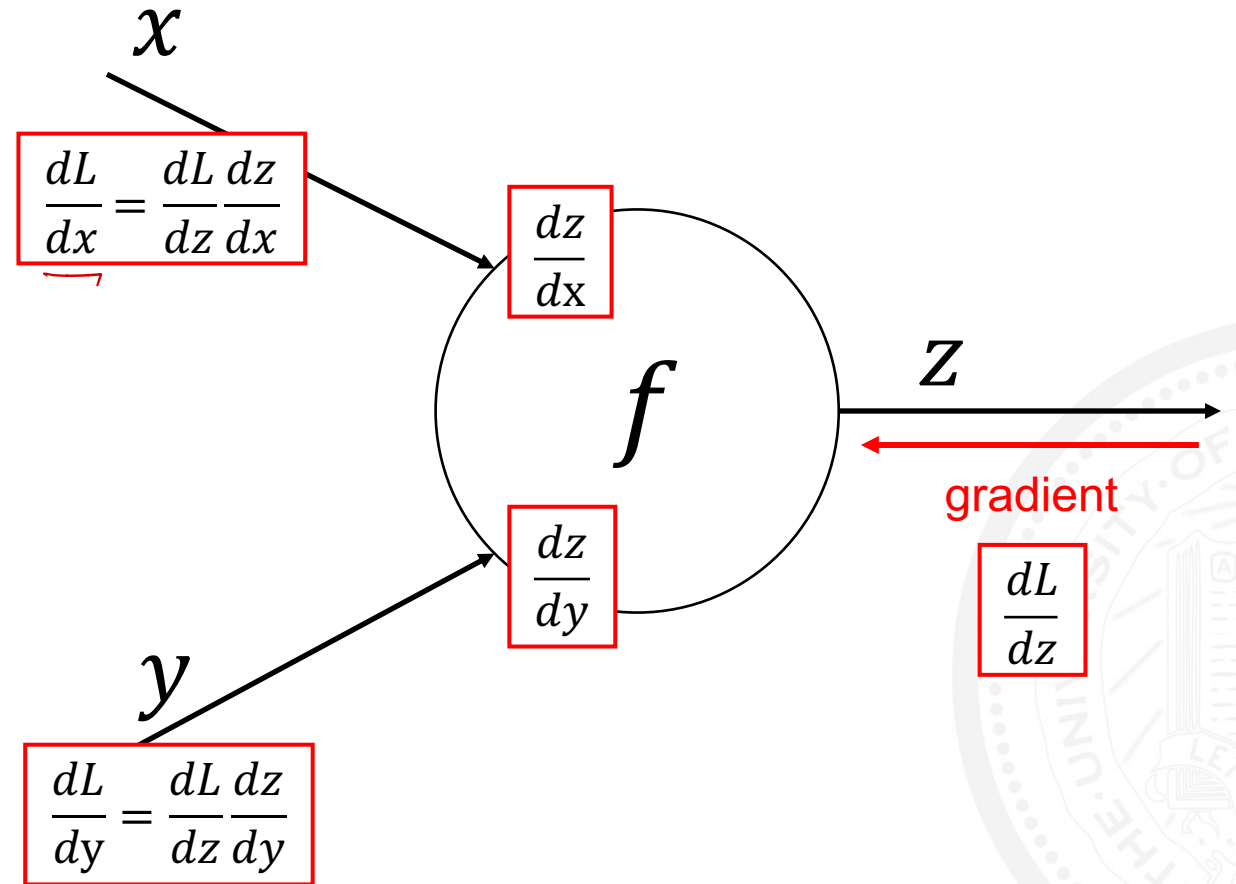
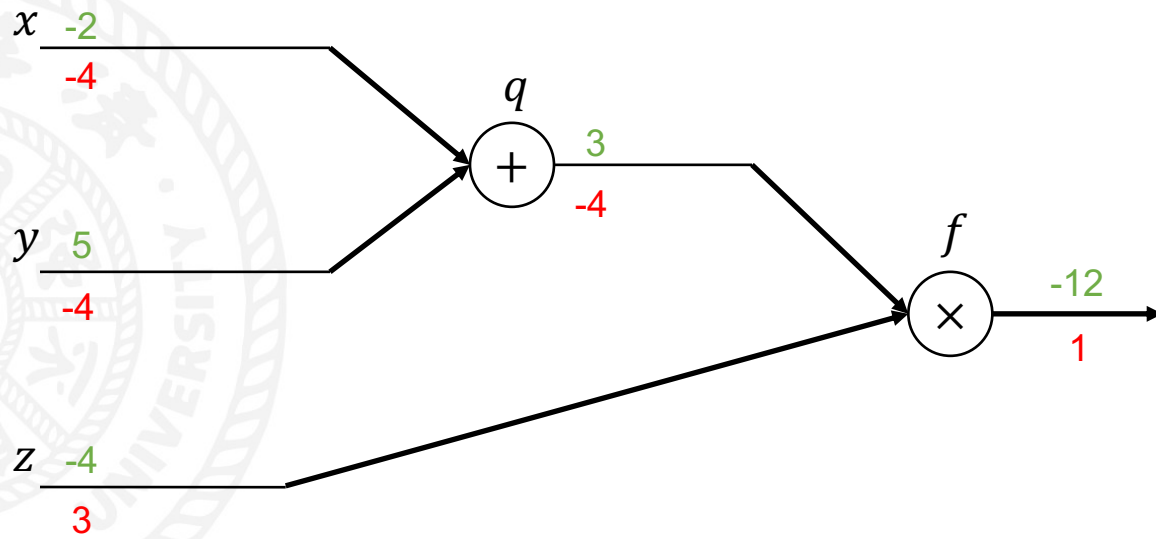
We want $\frac{df}{dx}$, $\frac{df}{dy}$, $\frac{df}{dz}$



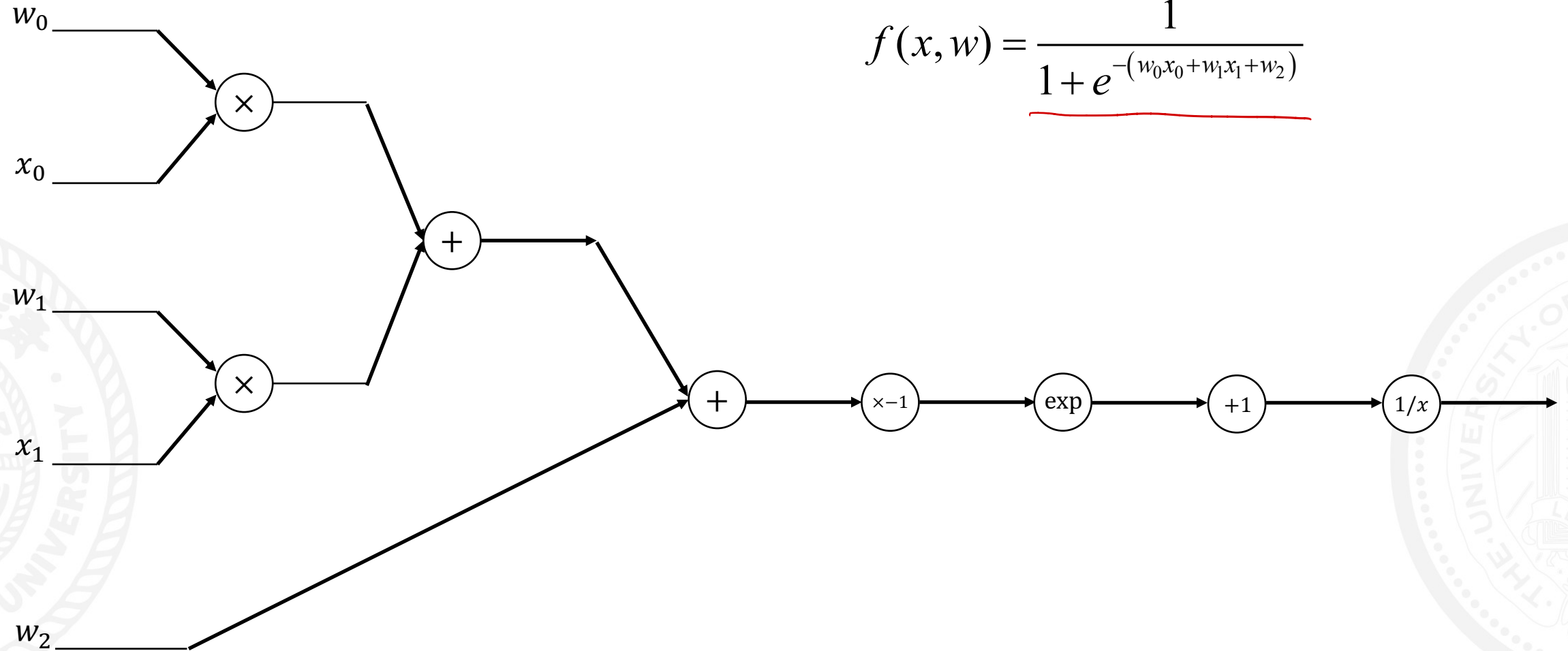
Backpropagation

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

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

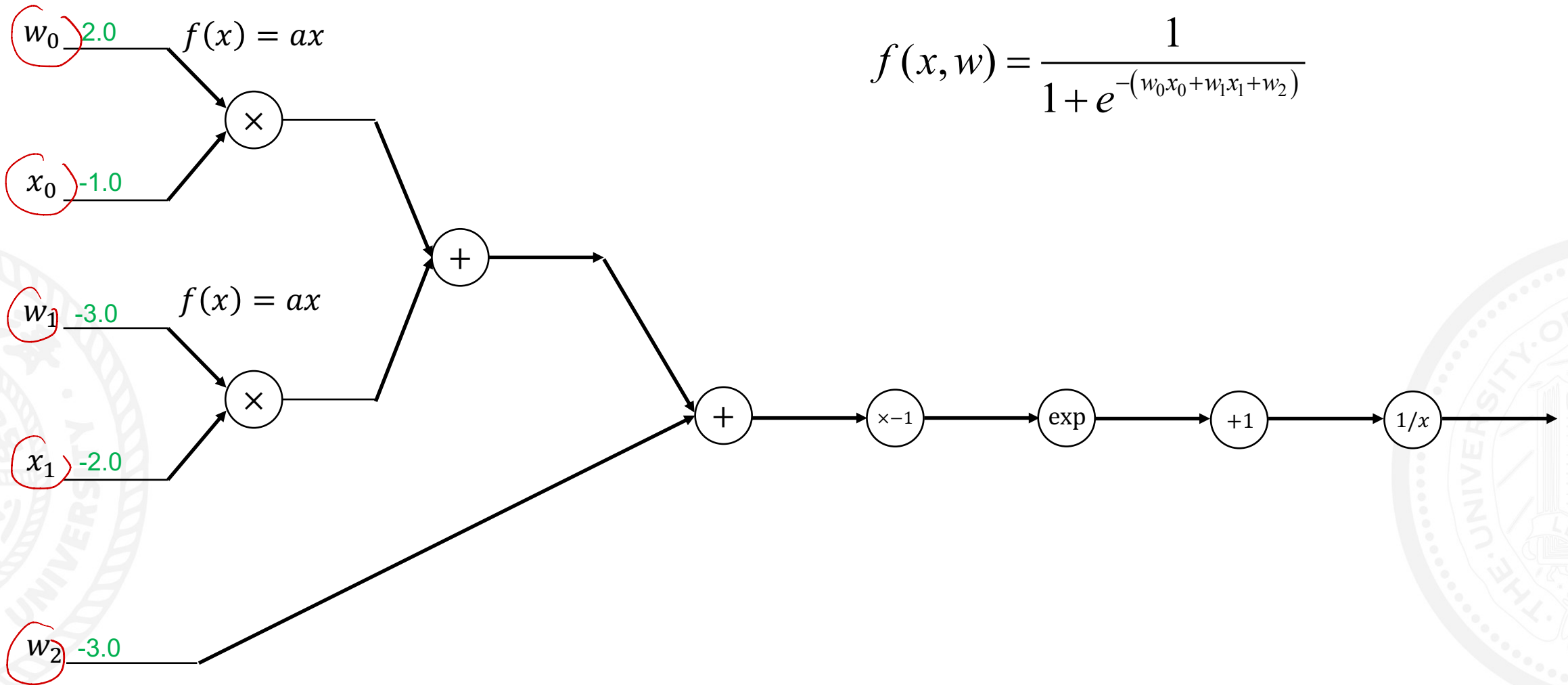


Backpropagation

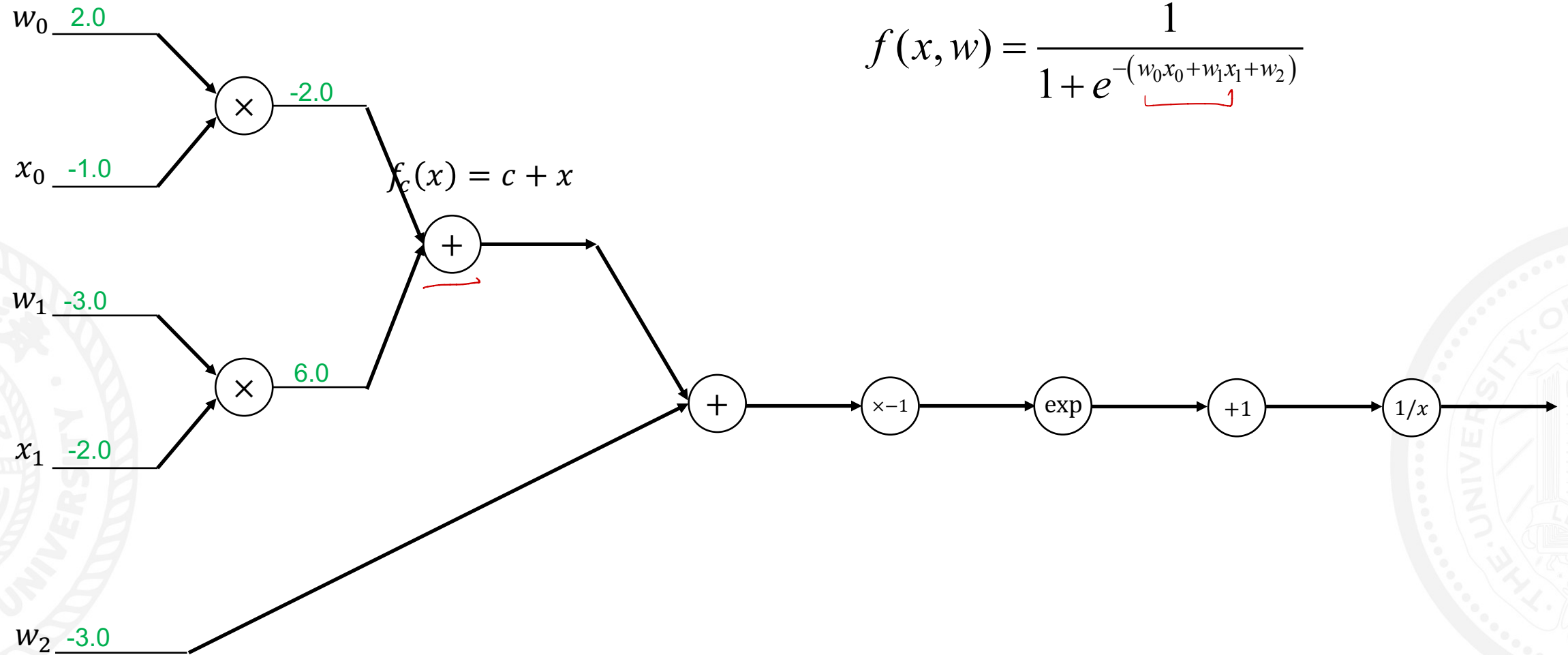


$$f(x, w) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

Backpropagation

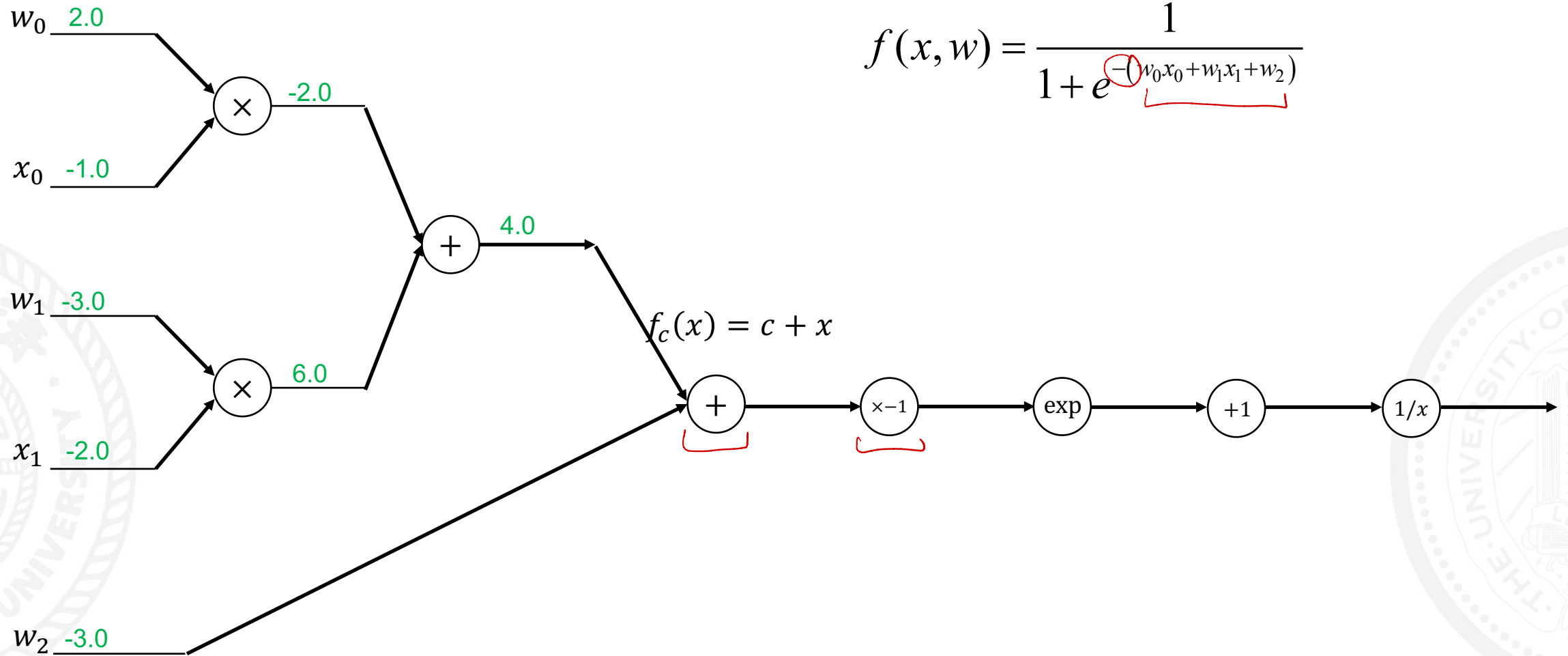


Backpropagation



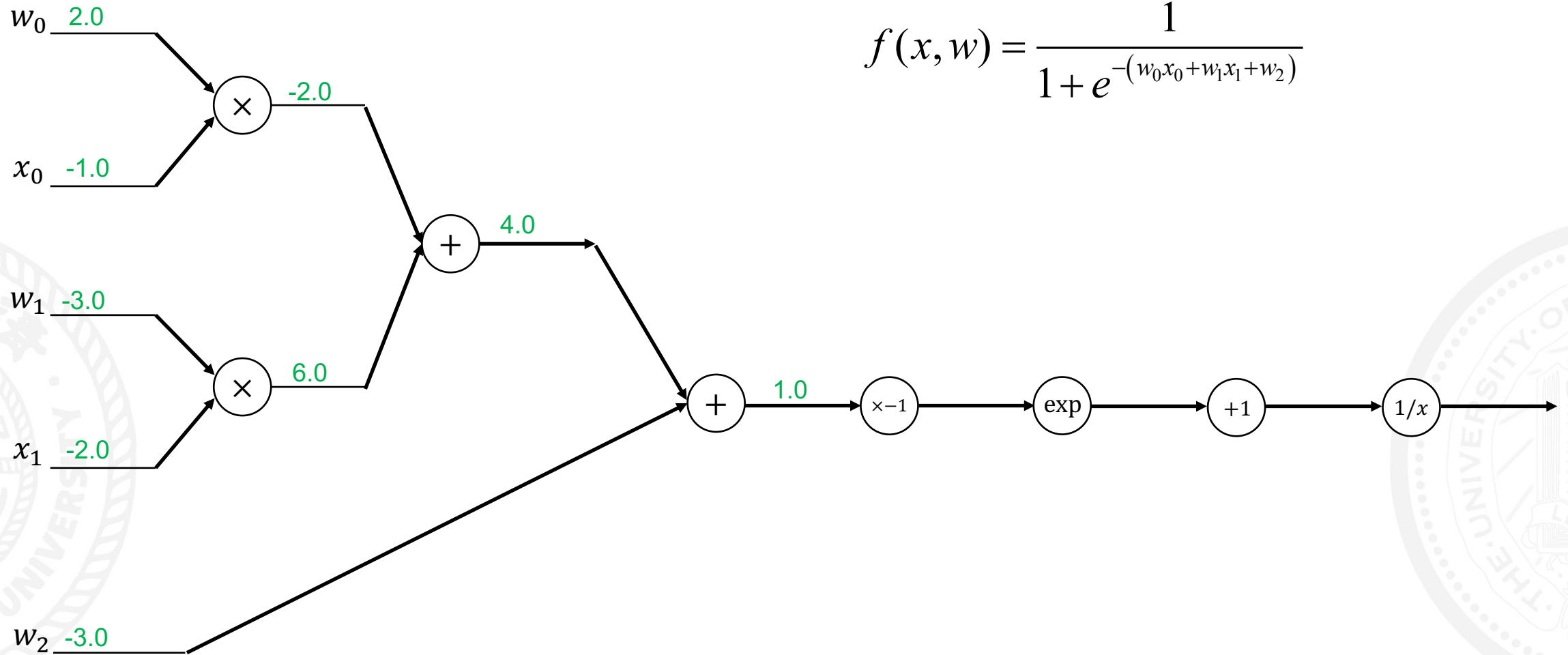
$$f(x, w) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

Backpropagation



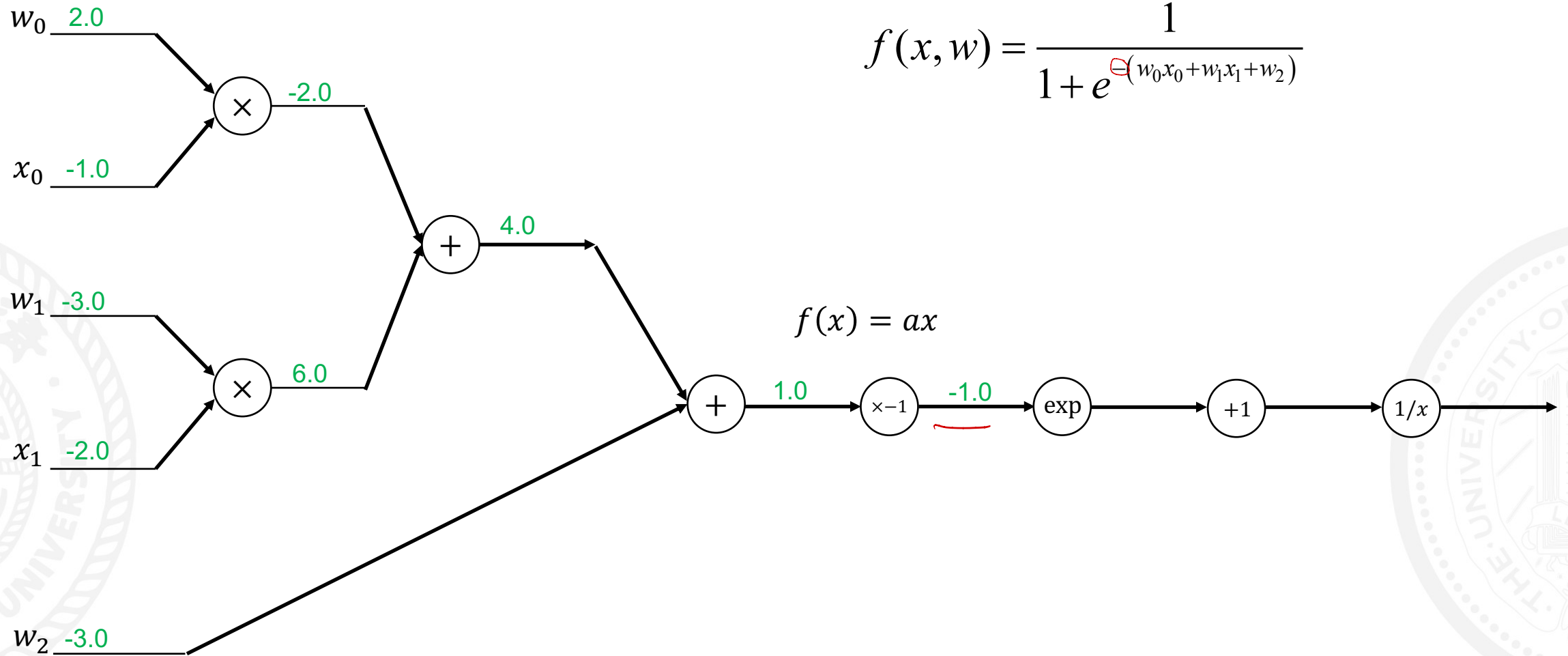
$$f(x, w) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

Backpropagation



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

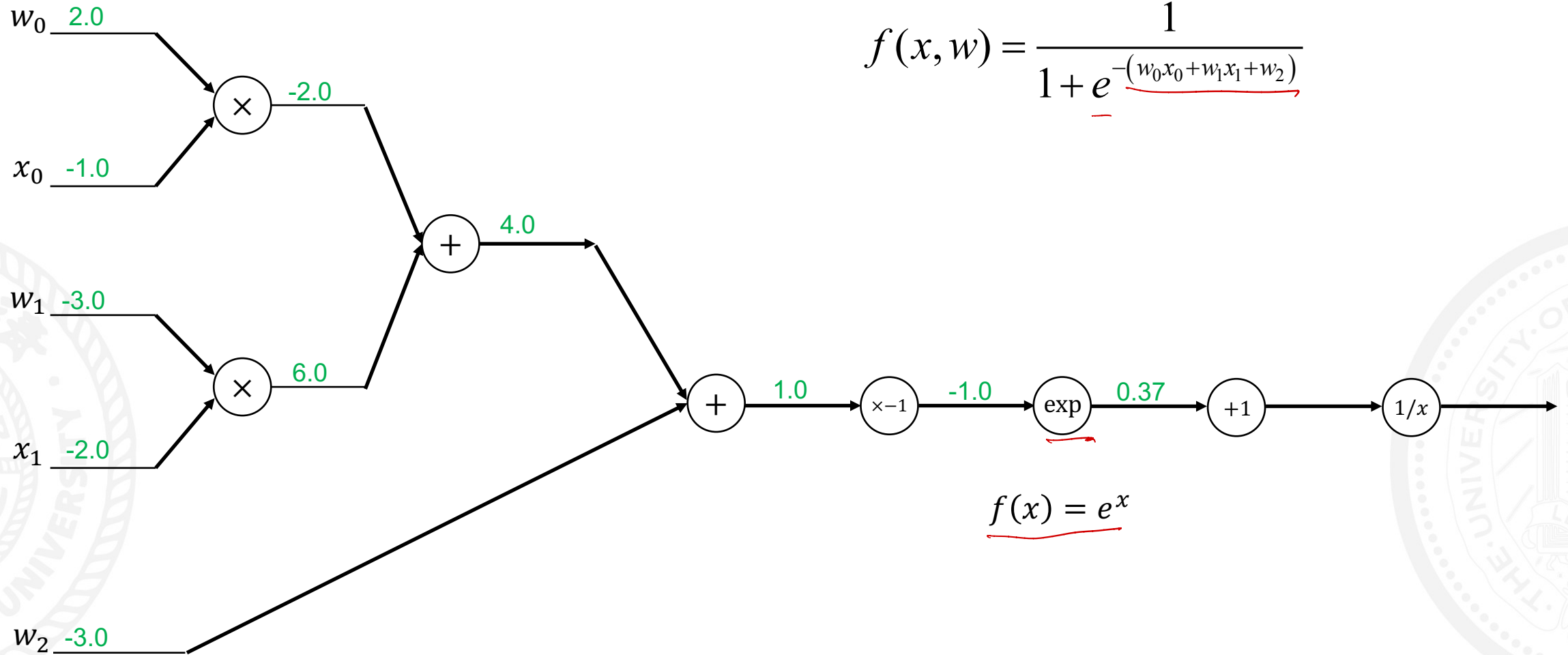
Backpropagation



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

$$f(x) = ax$$

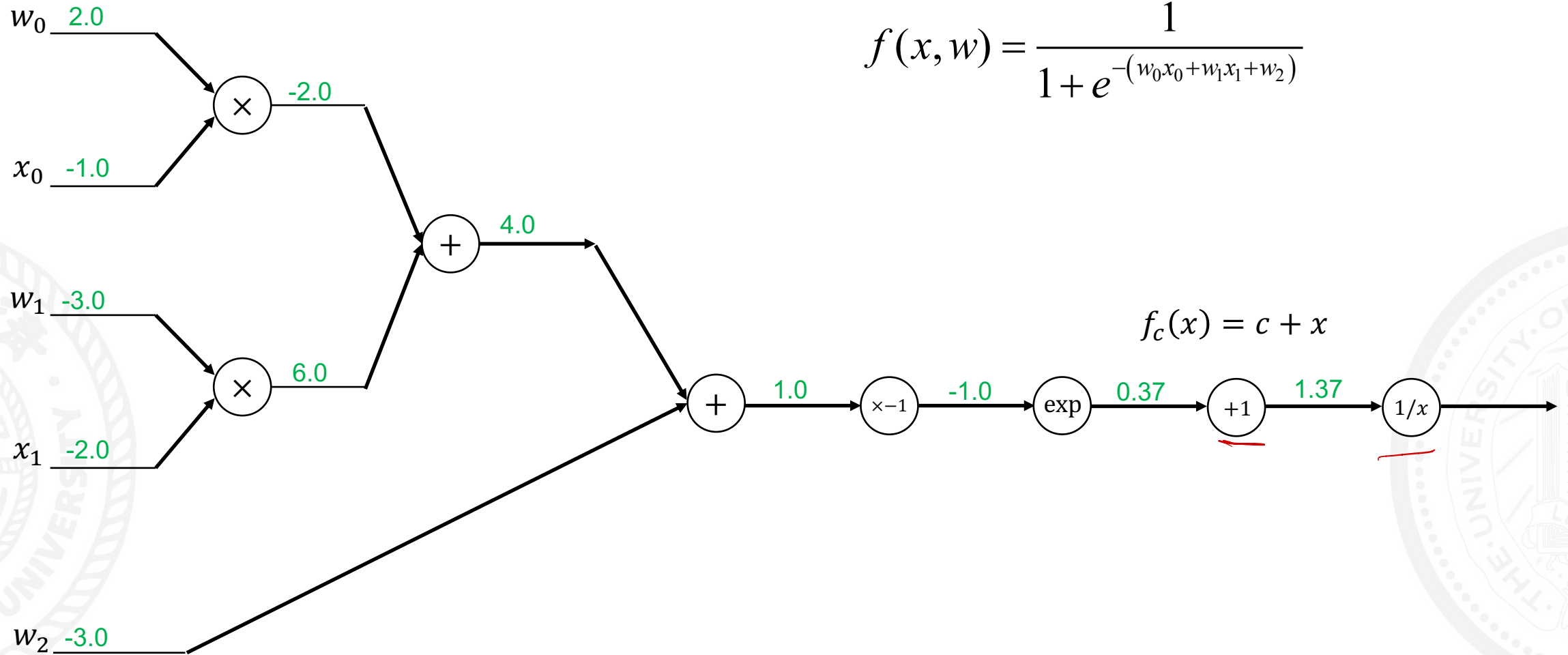
Backpropagation



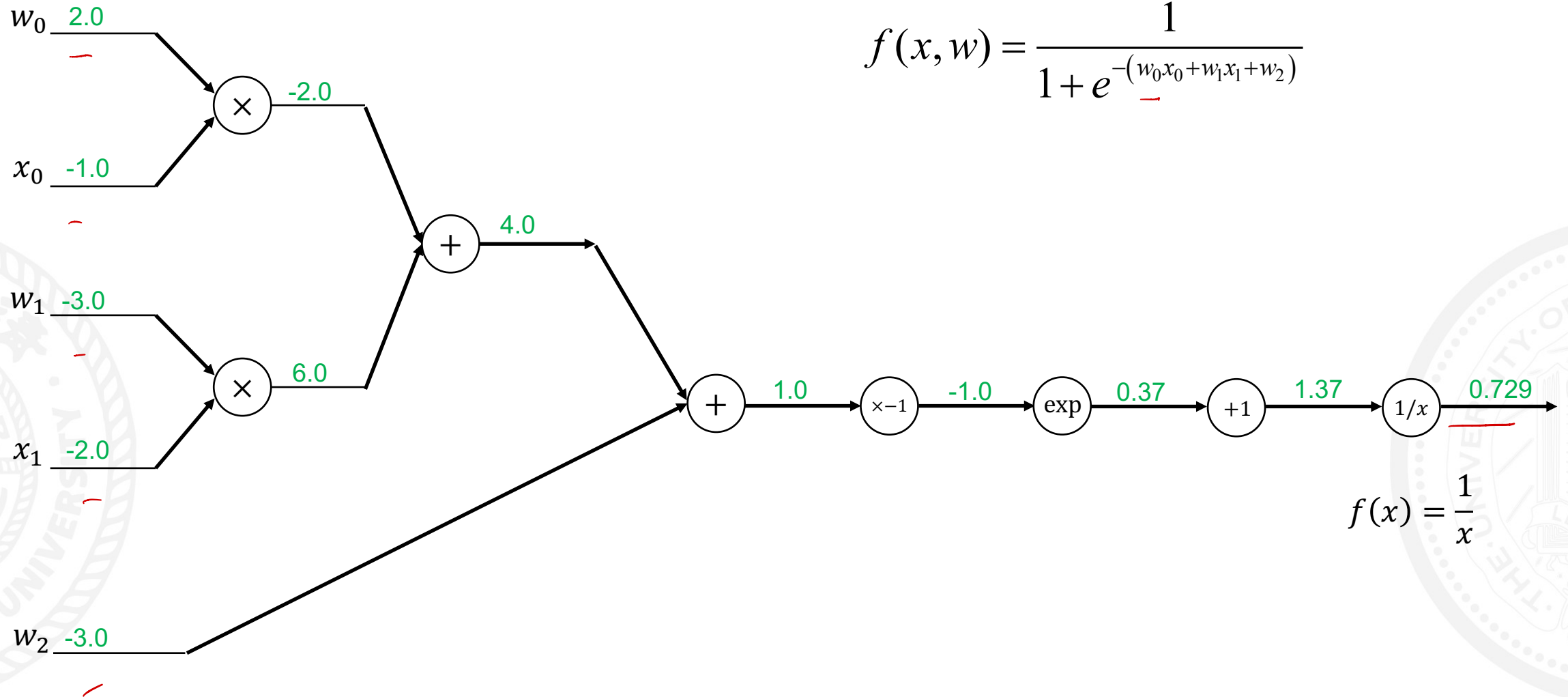
$$f(x, w) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

$$f(x) = e^x$$

Backpropagation



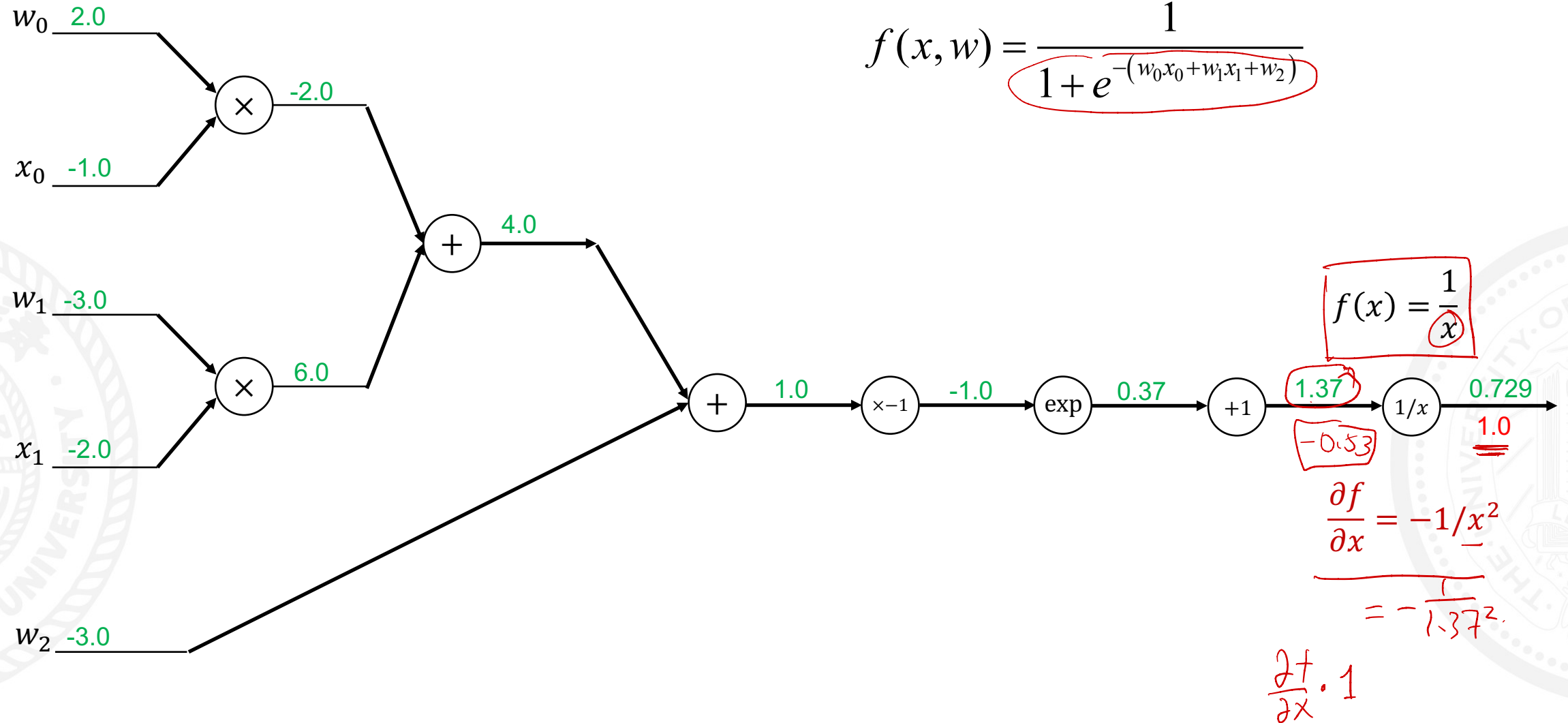
Backpropagation



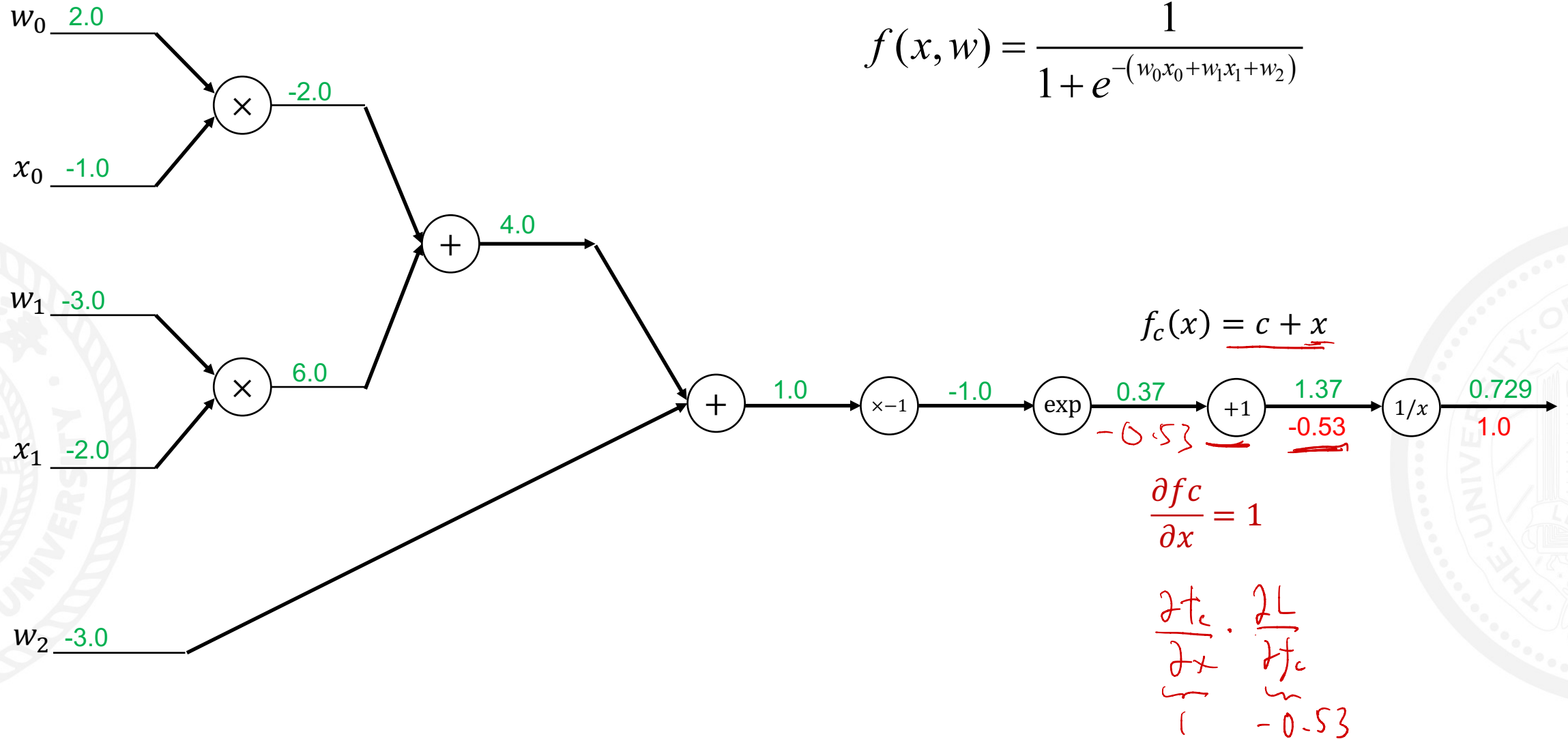
$$f(x, w) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

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

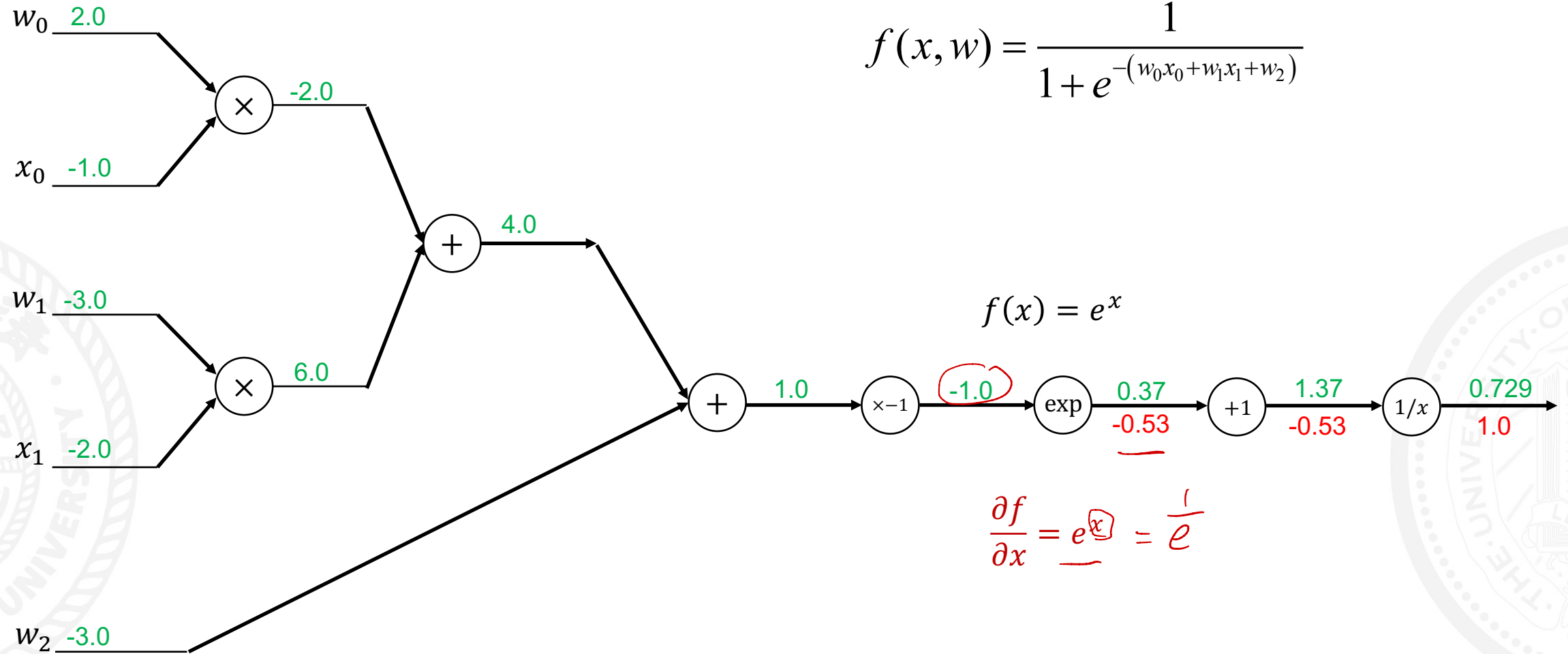
Backpropagation



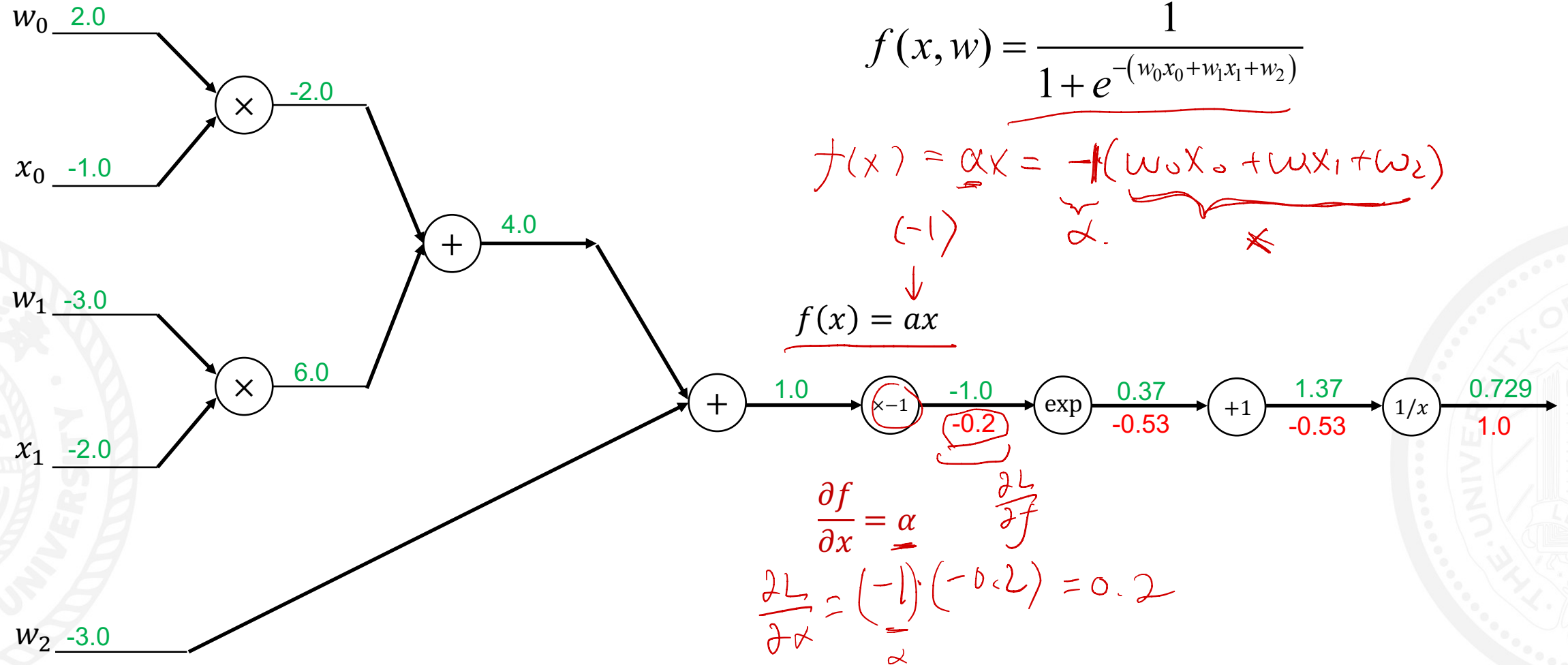
Backpropagation



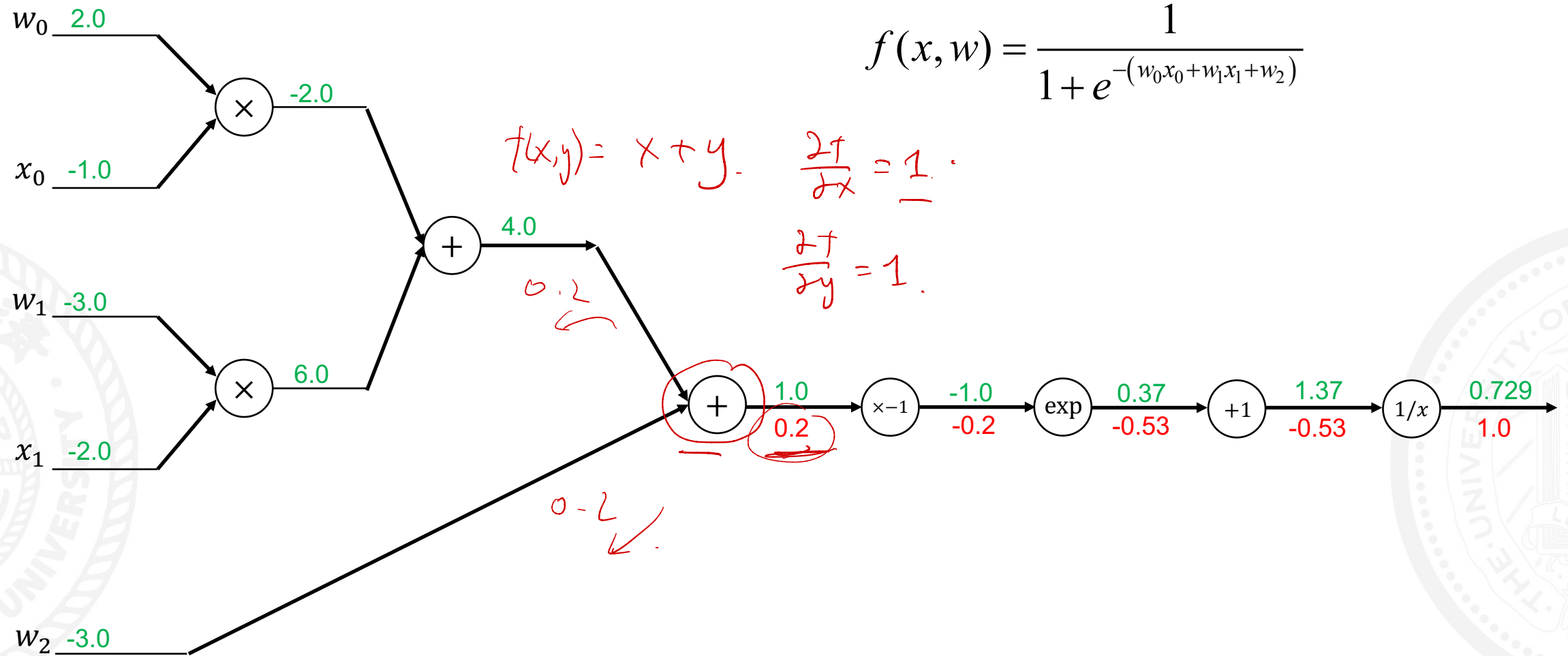
Backpropagation



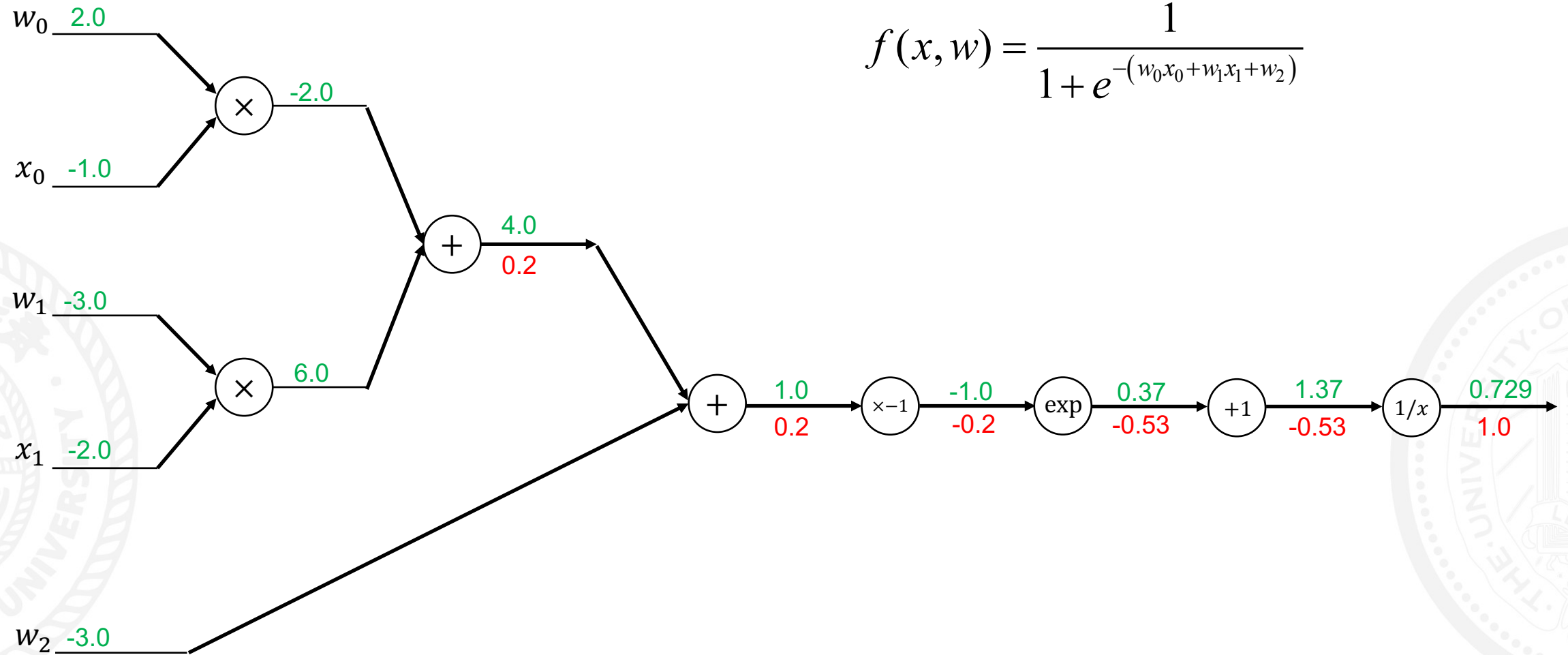
Backpropagation



Backpropagation

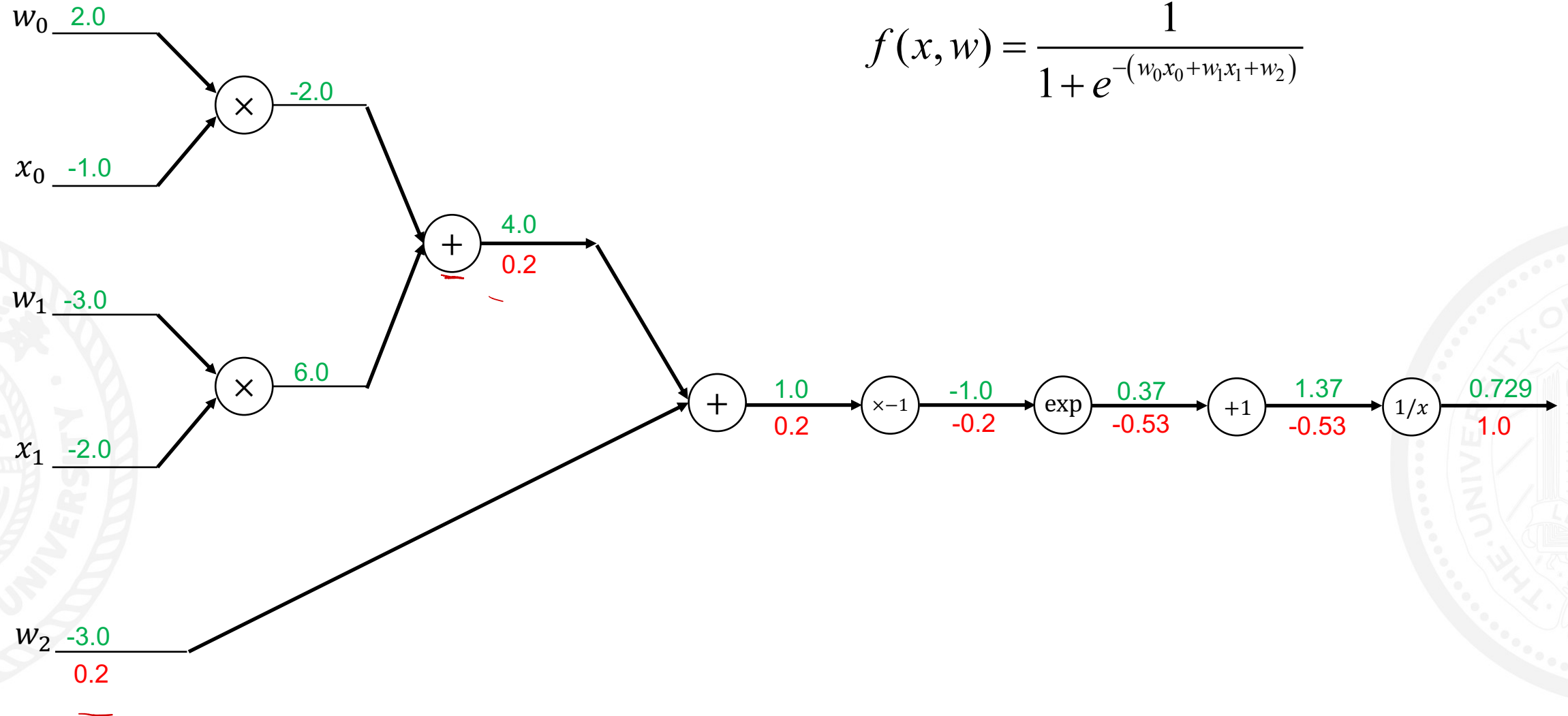


Backpropagation



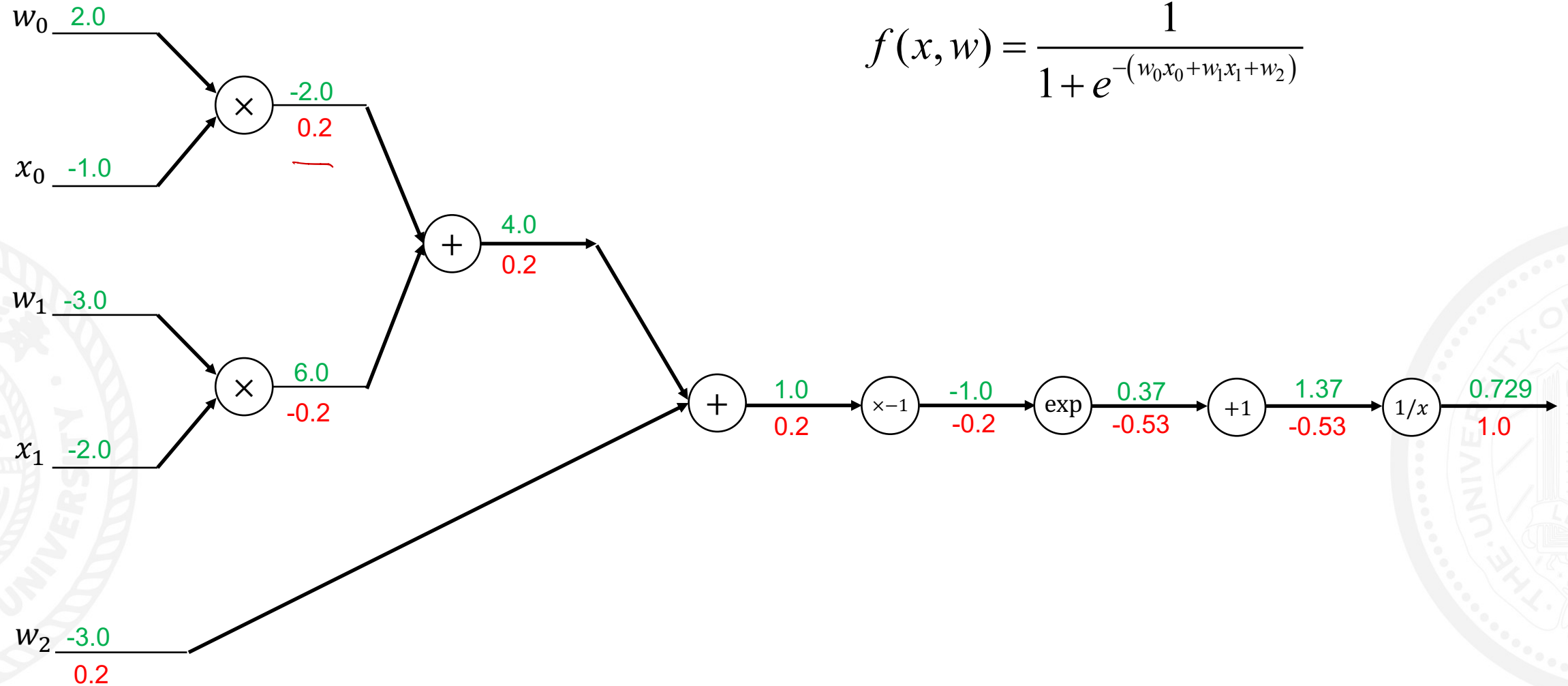
$$f(x, w) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

Backpropagation



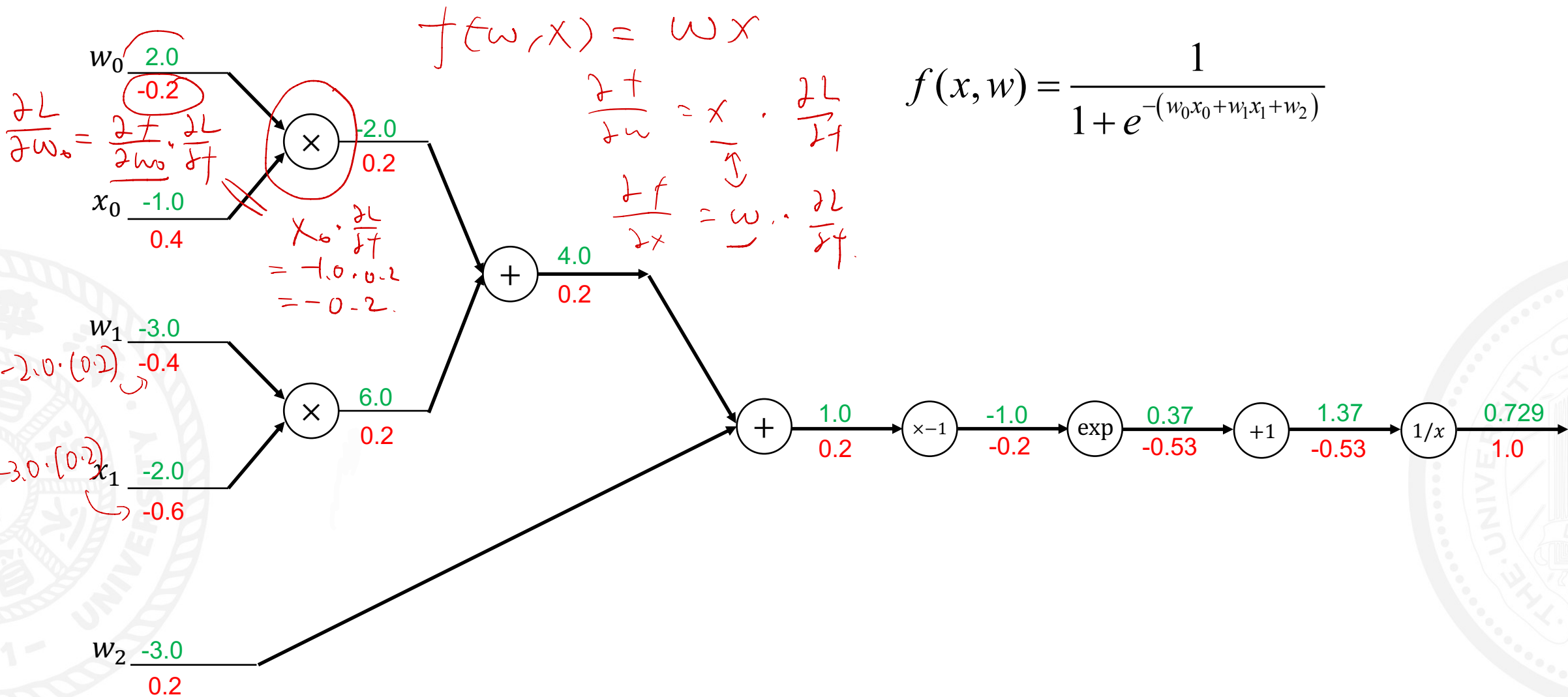
$$f(x, w) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

Backpropagation

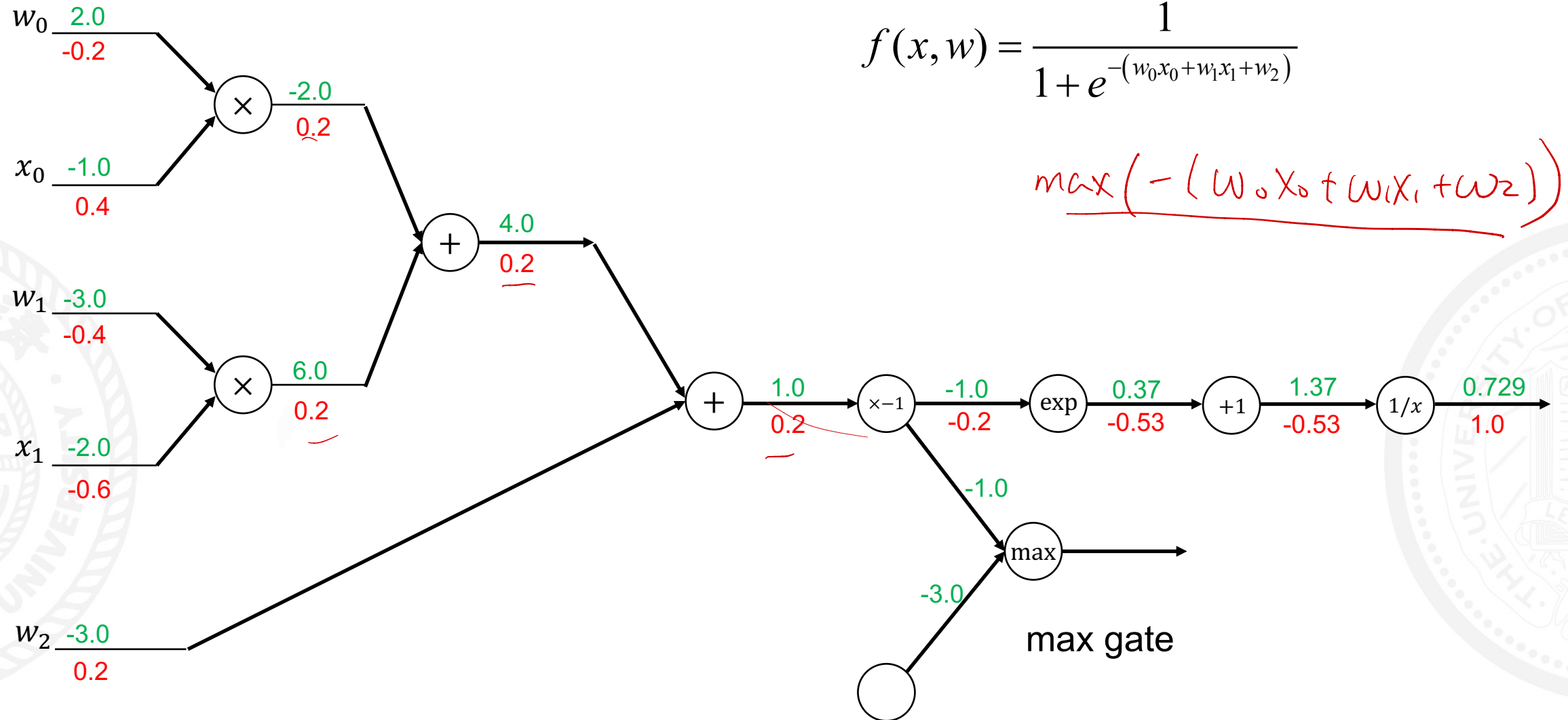


$$f(x, w) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

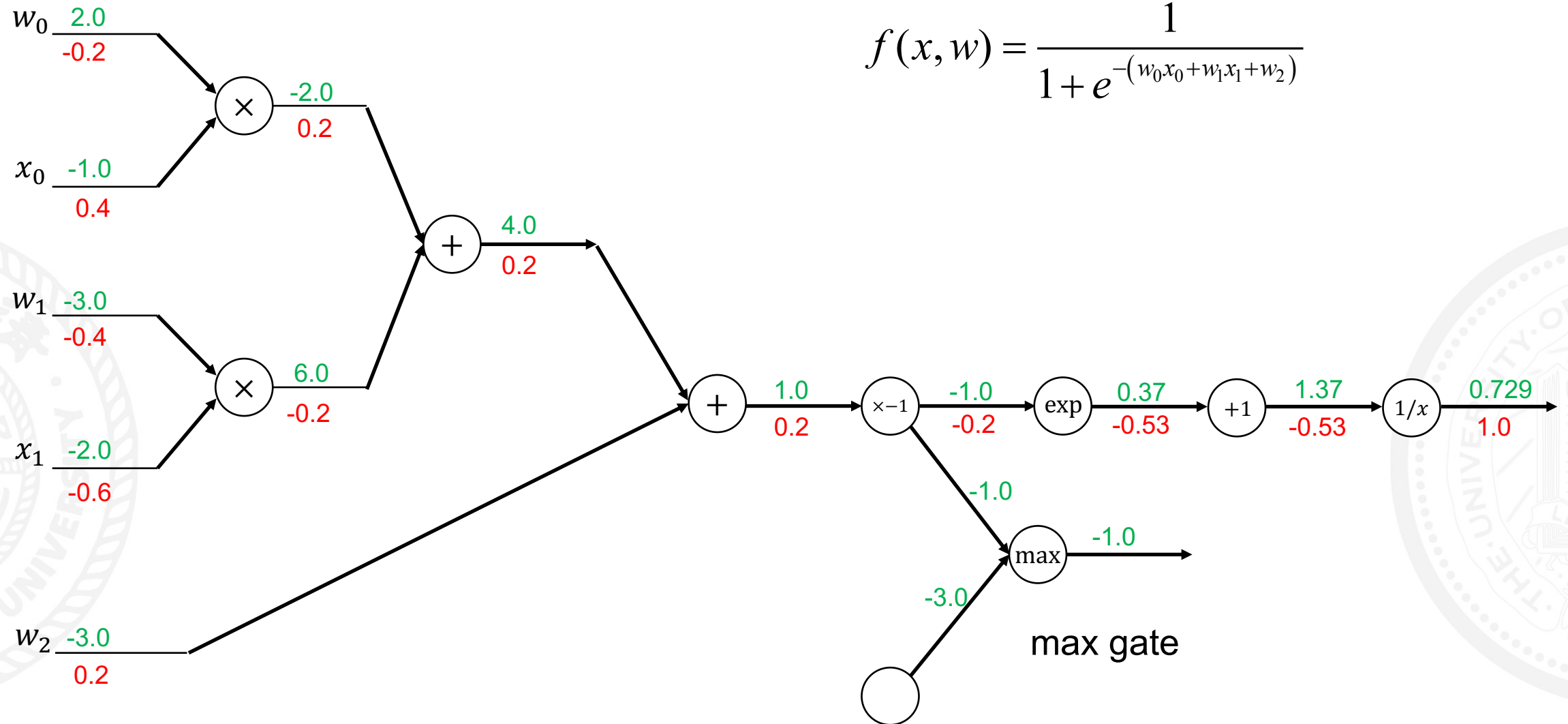
Backpropagation



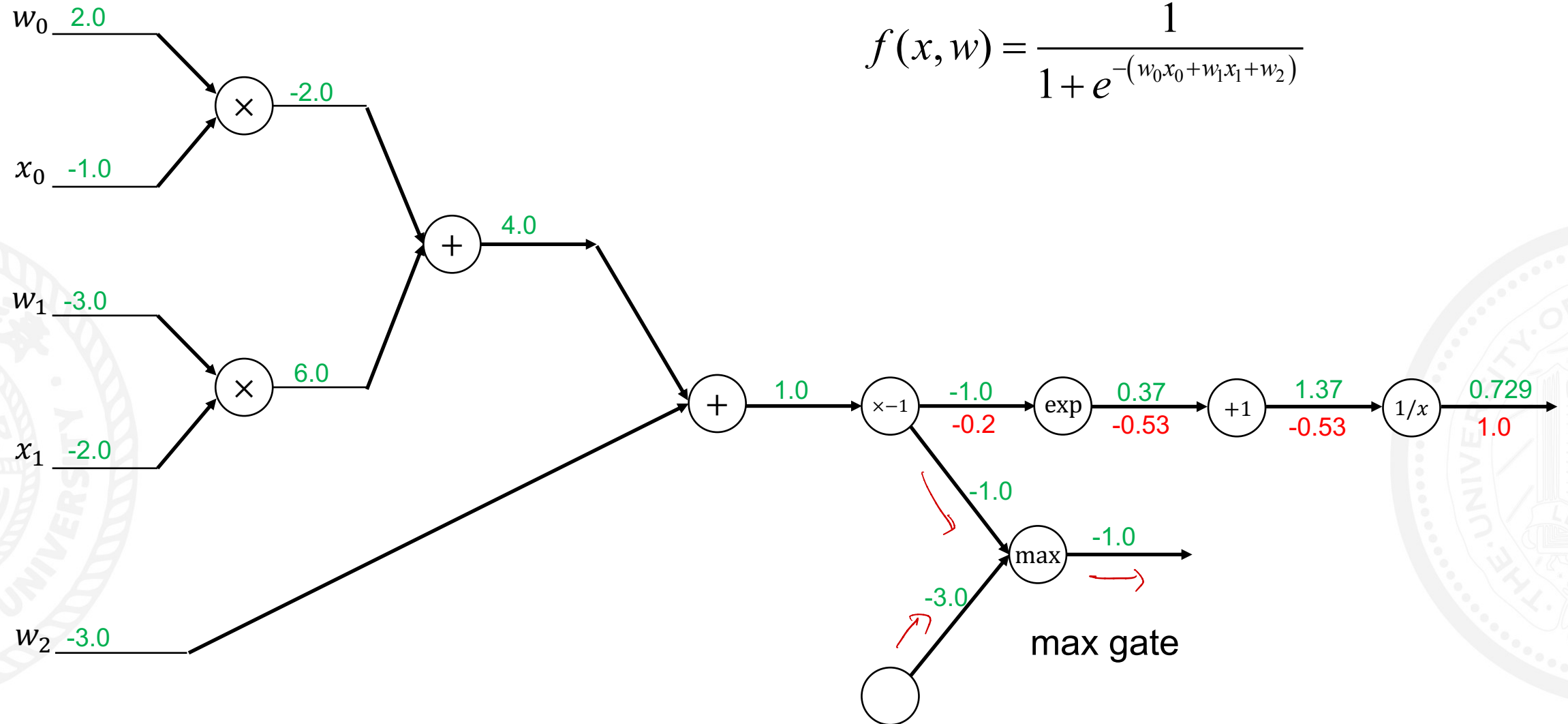
Backpropagation



Backpropagation



Backpropagation

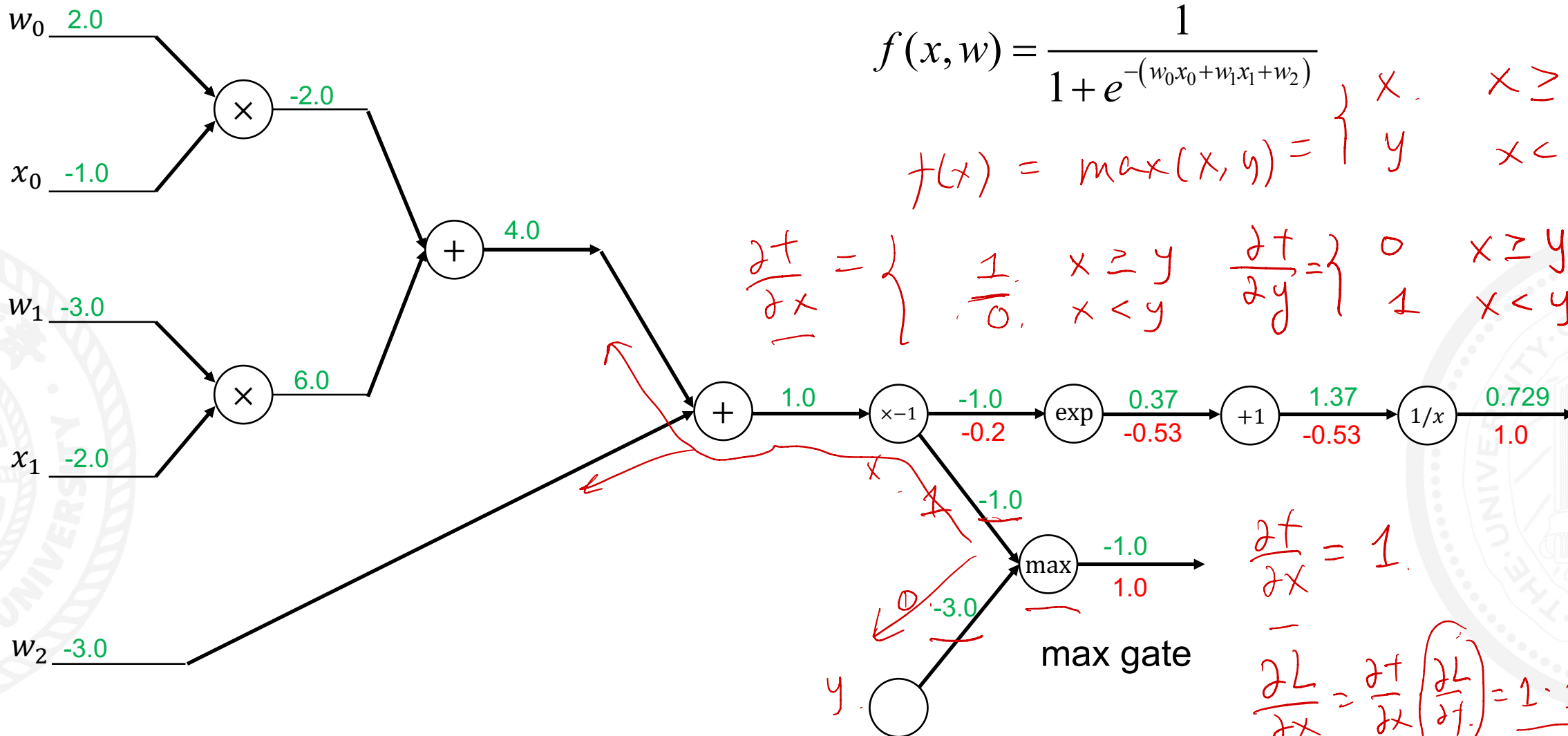


Backpropagation

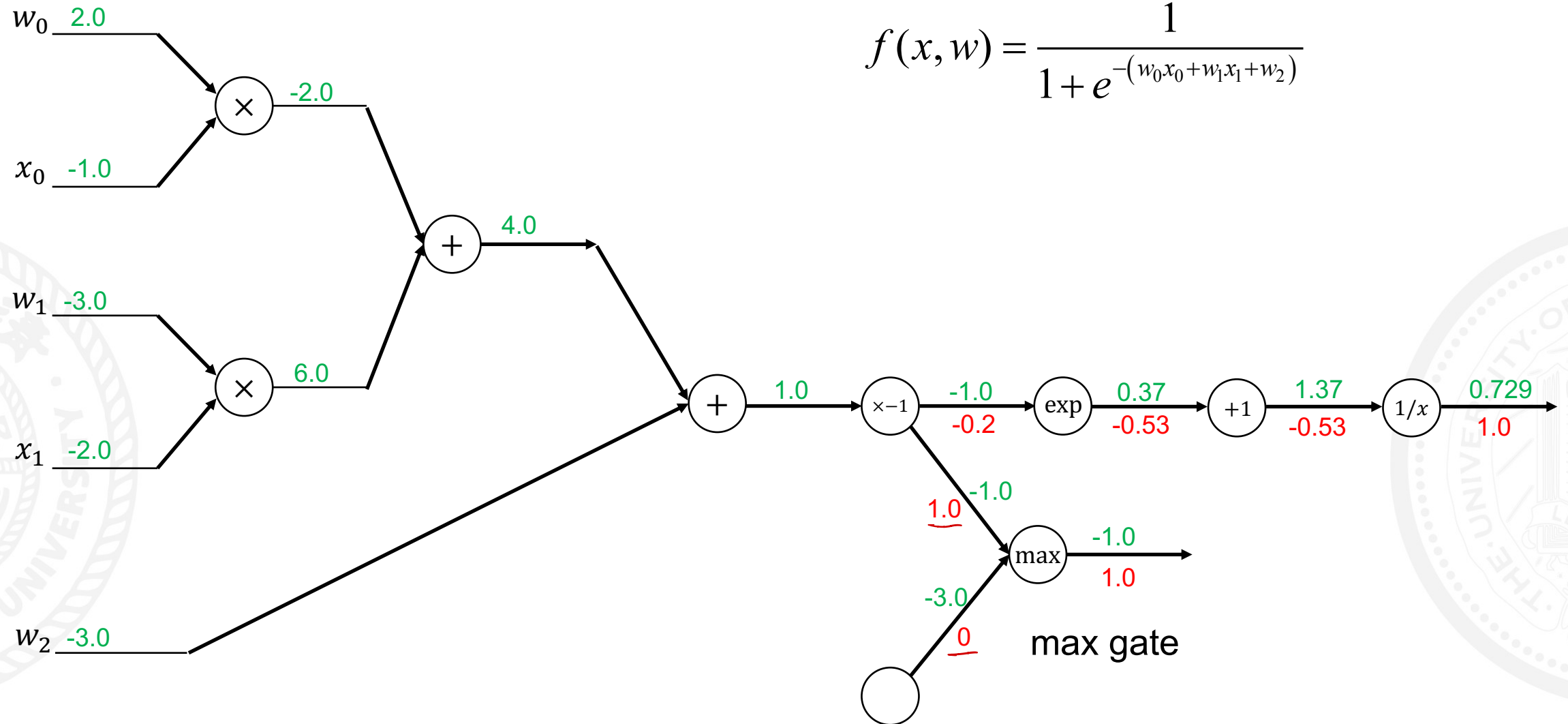


TBSI

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

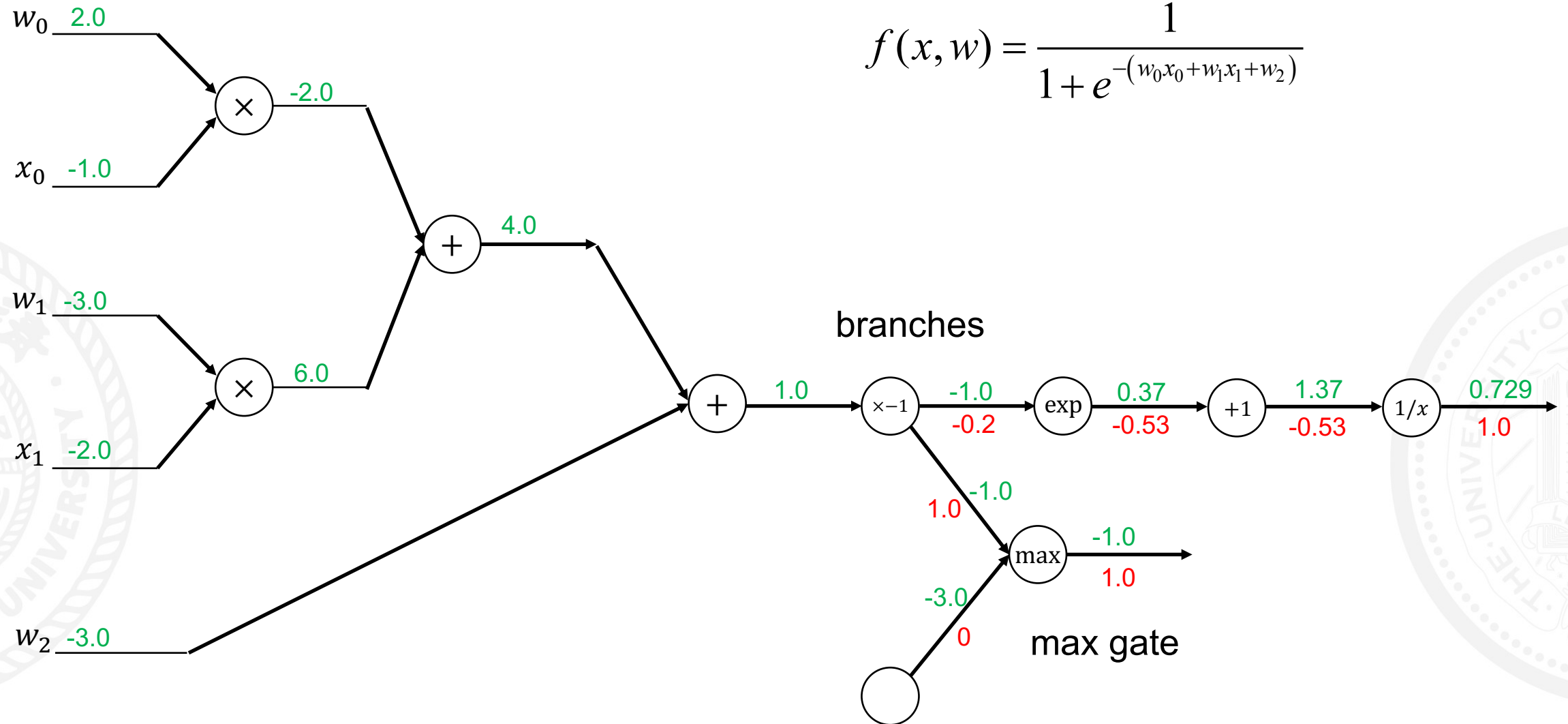


Backpropagation



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

Backpropagation

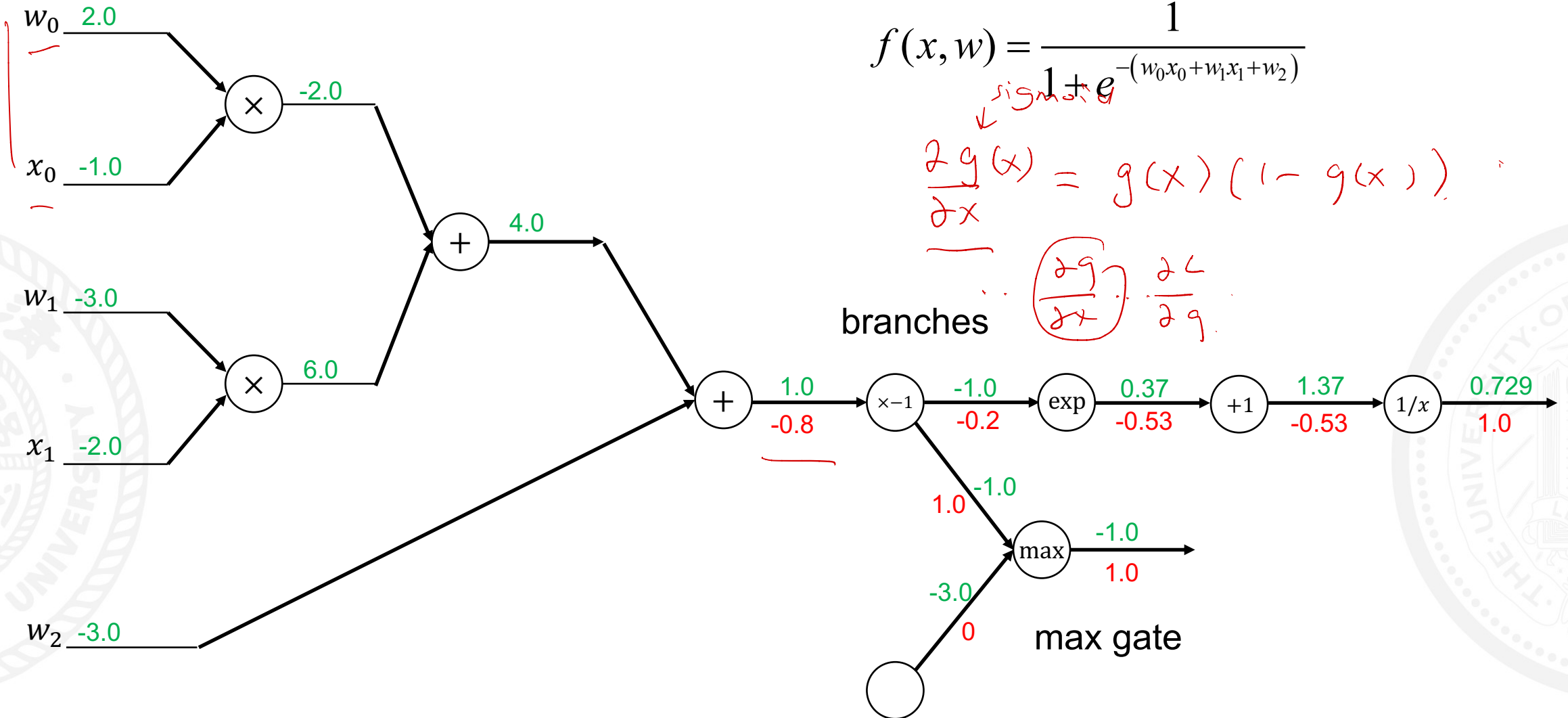


$$f(x, w) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

branches

max gate

Backpropagation

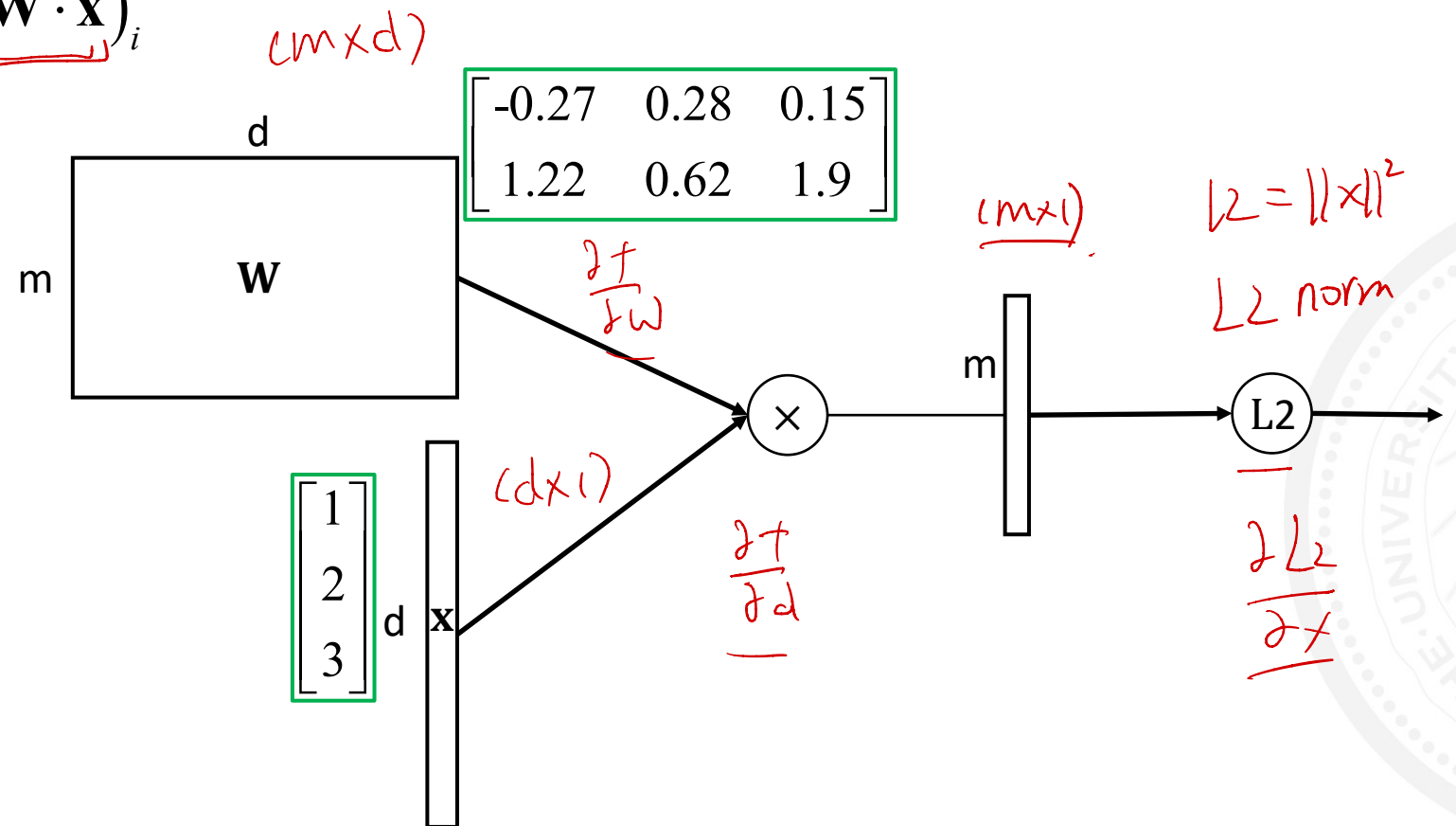


Backpropagation



Vectorized example

$$f(\mathbf{x}, \mathbf{W}) = \|\mathbf{W} \cdot \mathbf{x}\|^2 = \sum_{i=1}^n (\mathbf{W} \cdot \mathbf{x})_i^2$$

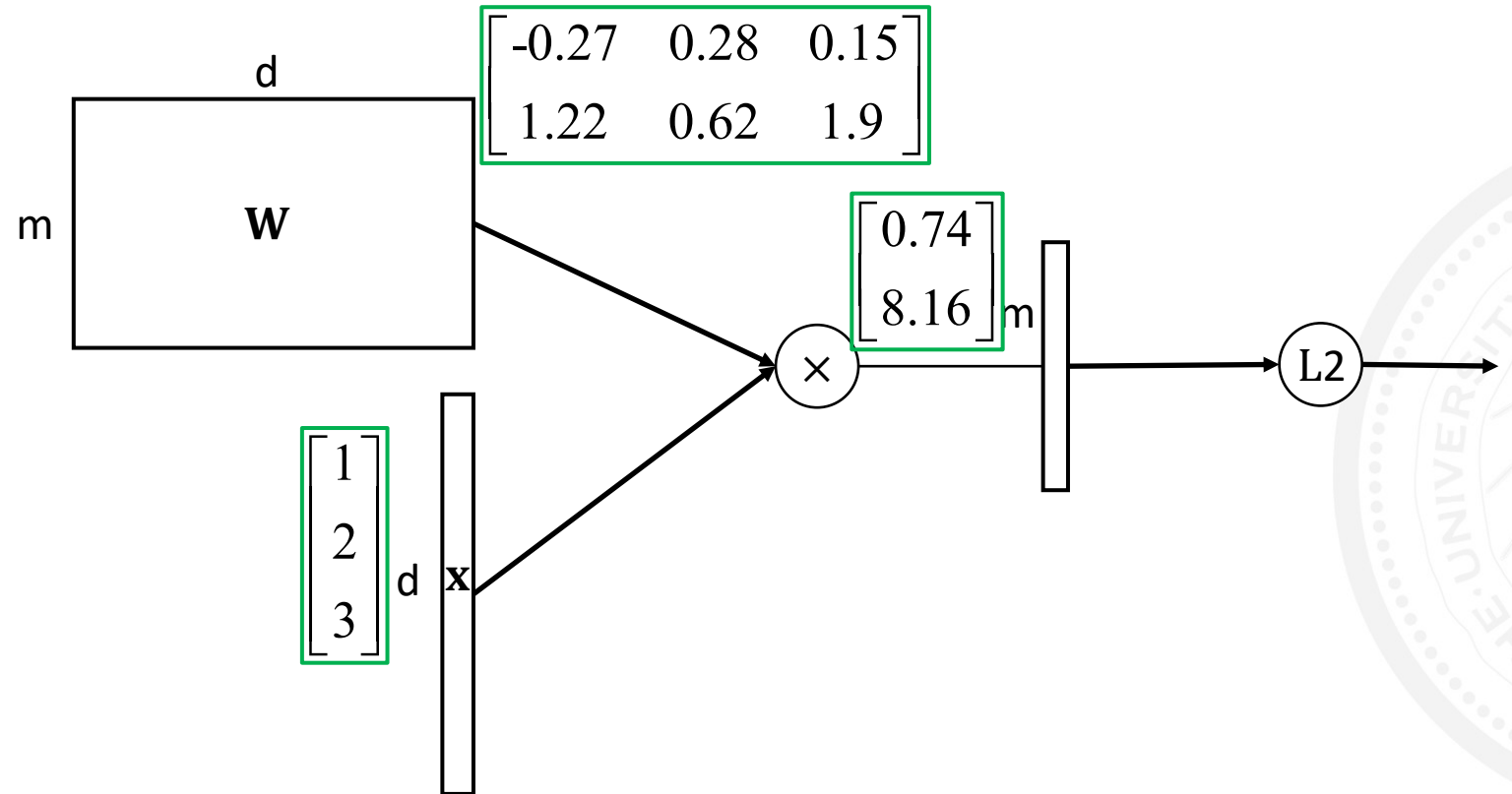


Backpropagation



Vectorized example

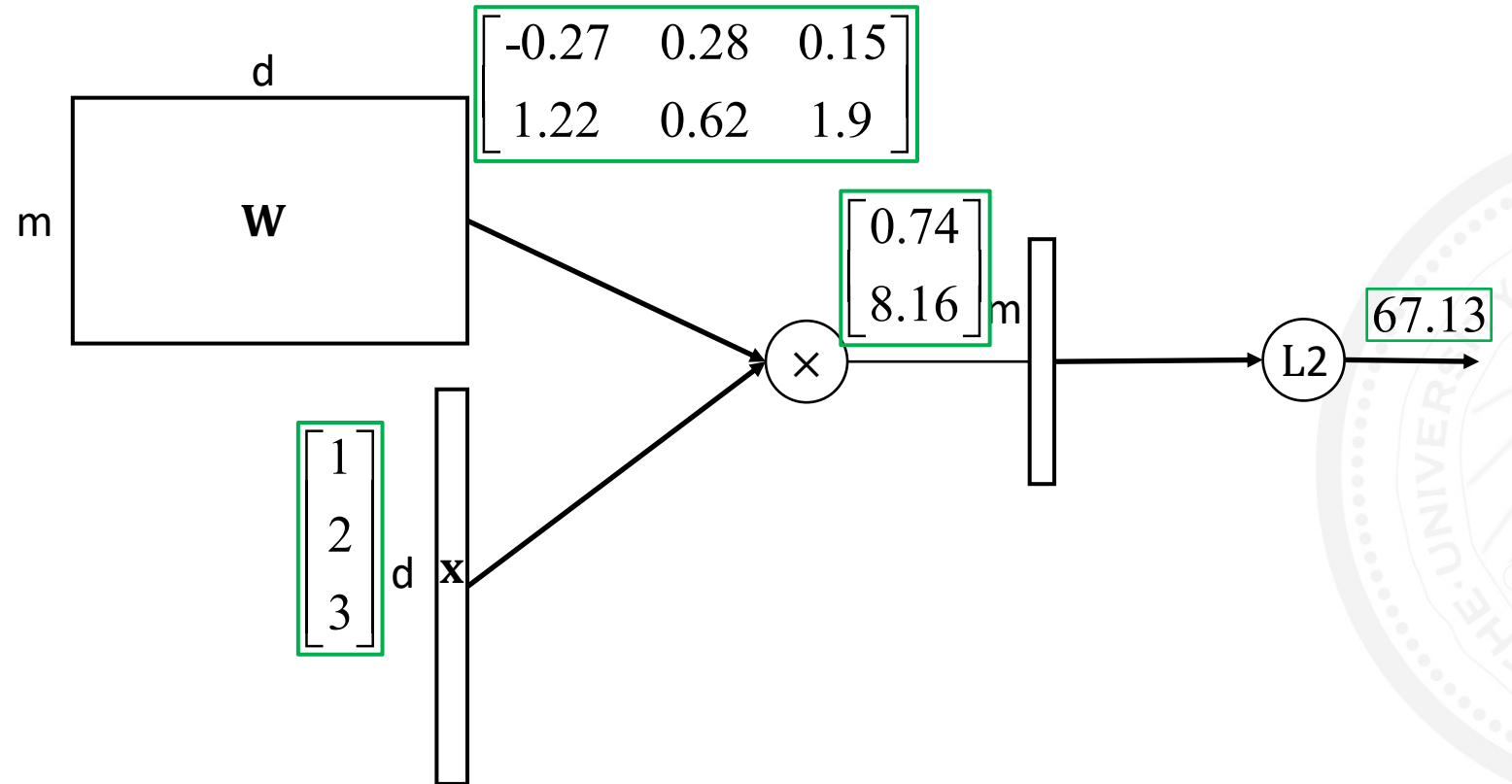
$$f(\mathbf{x}, \mathbf{W}) = \|\mathbf{W} \cdot \mathbf{x}\|^2 = \sum_{i=1}^n (\mathbf{W} \cdot \mathbf{x})_i^2$$



Backpropagation

Vectorized example

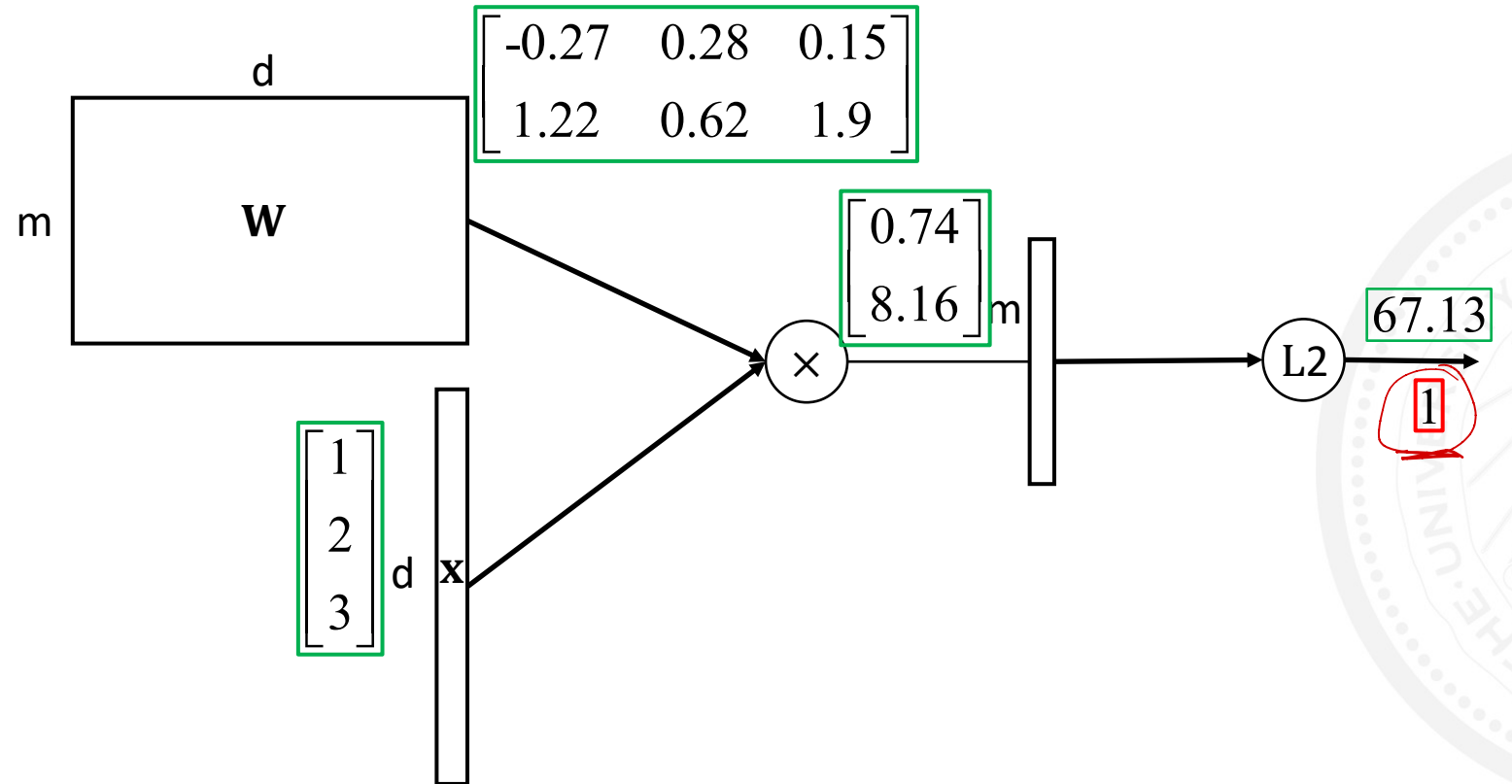
$$f(\mathbf{x}, \mathbf{W}) = \|\mathbf{W} \cdot \mathbf{x}\|^2 = \sum_{i=1}^n (\mathbf{W} \cdot \mathbf{x})_i^2$$



Backpropagation

Vectorized example

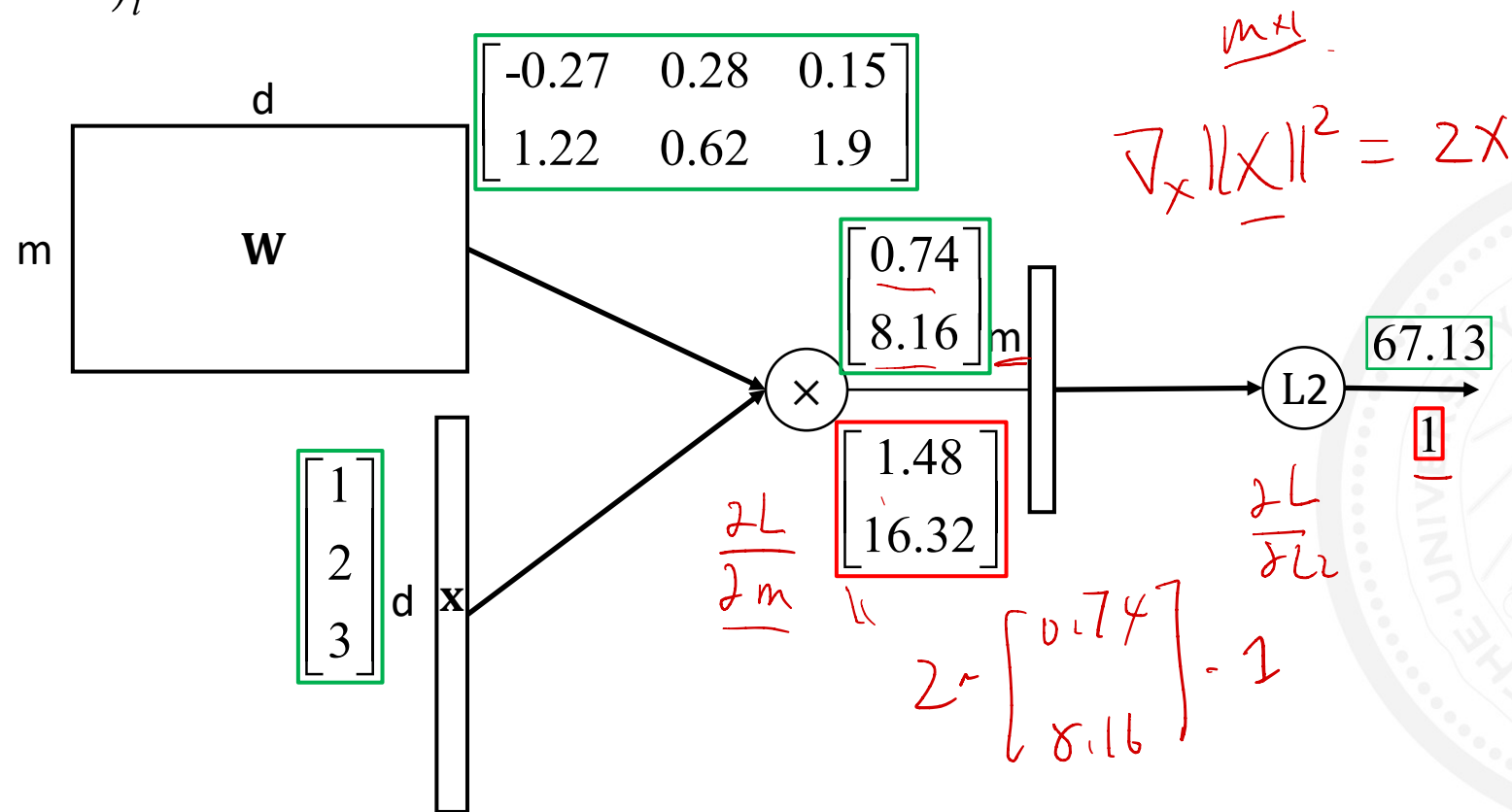
$$f(\mathbf{x}, \mathbf{W}) = \|\mathbf{W} \cdot \mathbf{x}\|^2 = \sum_{i=1}^n (\mathbf{W} \cdot \mathbf{x})_i^2$$



Backpropagation

Vectorized example

$$f(\mathbf{x}, \mathbf{W}) = \|\mathbf{W} \cdot \mathbf{x}\|^2 = \sum_{i=1}^n (\mathbf{W} \cdot \mathbf{x})_i^2$$



Backpropagation



Vectorized example

$$f(\mathbf{x}, \mathbf{W}) = \|\mathbf{W} \cdot \mathbf{x}\|^2 = \sum_{i=1}^n (\mathbf{W} \cdot \mathbf{x})_i^2$$

① $a(\mathbf{x}, \mathbf{W}) = \mathbf{W}^T \mathbf{x} = \mathbf{W} \cdot \mathbf{x}$ $\left\{ \begin{array}{l} \nabla_{\mathbf{W}} a(\mathbf{x}, \mathbf{W}) = \mathbf{x}^T \\ \nabla_{\mathbf{x}} a(\mathbf{x}, \mathbf{W}) = \mathbf{W} \end{array} \right\}$ local gradient

② $f(\mathbf{a}) = \|\mathbf{a}\|^2 = \mathbf{a}^T \mathbf{a}$ $\left\{ \begin{array}{l} \nabla_{\mathbf{a}} f(\mathbf{a}) \\ = 2\mathbf{a} \end{array} \right\}$ local gradient

$(m \times 1)$ $(1 \times d)$ $\nabla_{\mathbf{W}} a(\mathbf{W}, \mathbf{x}) = \mathbf{x}^T$

$$\begin{bmatrix} 1.48 \\ 16.32 \end{bmatrix} \times \begin{bmatrix} 1 & 2 & 3 \end{bmatrix} = \begin{bmatrix} 1.48 & 2.96 & 4.44 \\ 16.32 & 32.64 & 48.96 \end{bmatrix}$$

$$\nabla_{\mathbf{a}} f(\mathbf{a}) \cdot \nabla_{\mathbf{W}} a(\mathbf{x}, \mathbf{W}) = \nabla_{\mathbf{W}} f(\mathbf{x}, \mathbf{W})$$

\mathbf{x}^T

