

# Planning, Learning and Decision Making

Lecture 15. Learning from examples - Neural approaches

# Approaches to learning

- “Symbolic” approach → Learn rules (DT)
- Probabilistic approach → Learn probabilities (LR, NB)
- Similarity-based approach → Learn by similarity (kNN, SVM)



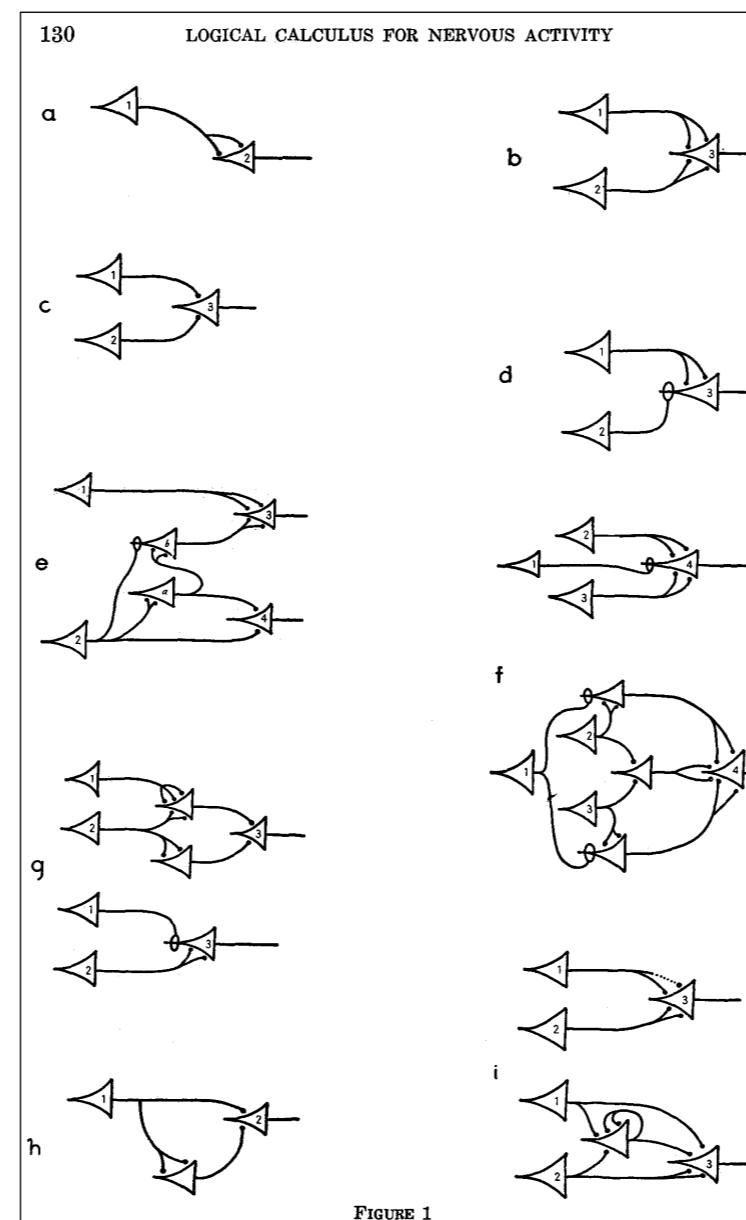
Neural approach

# Neural approach

- Inspired by how our brain works
- Idea has been around since 1940s (age of “cybernetics”)

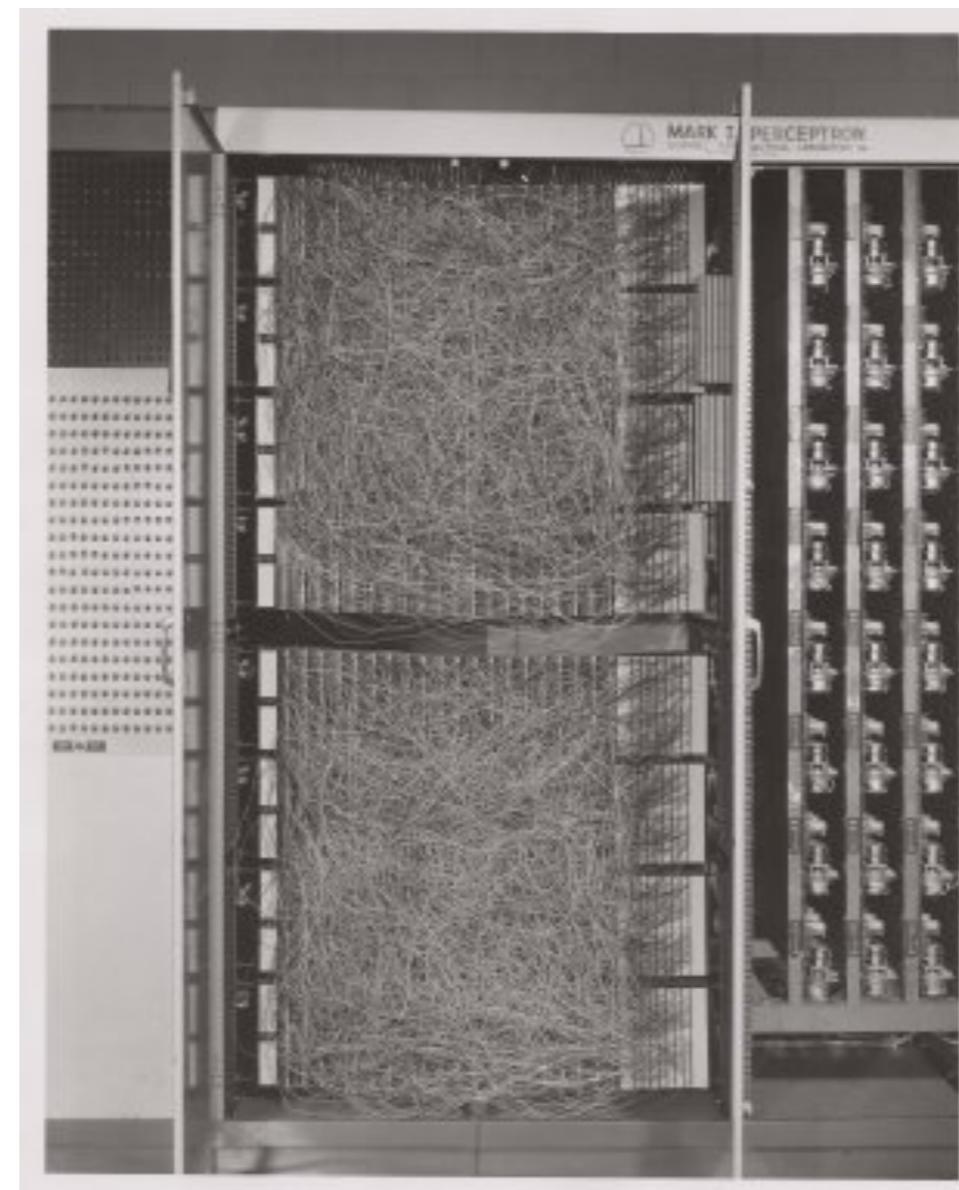
# Linear Threshold Unit

- McCulloch & Pitts (1943)



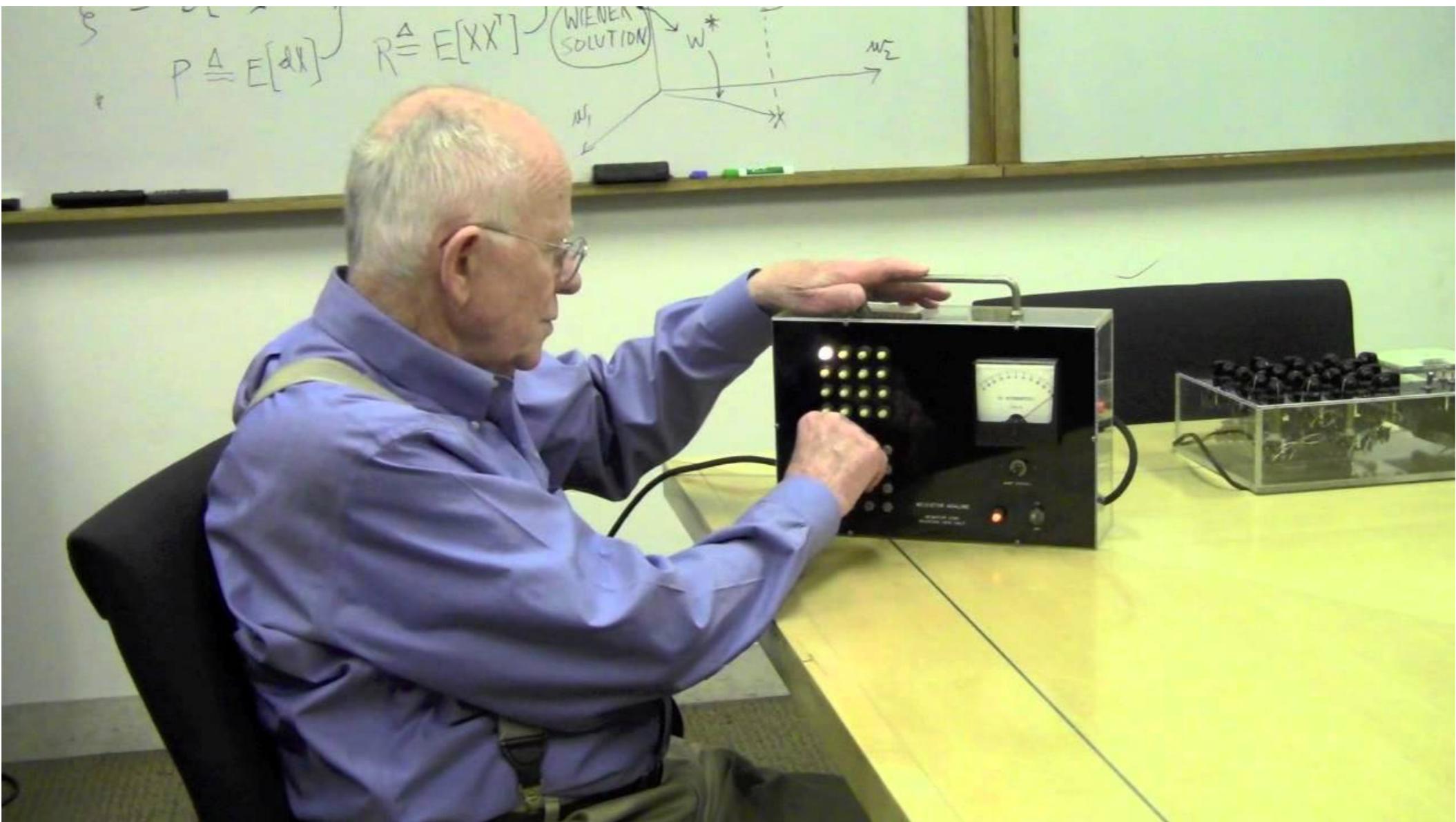
# Perceptron

- Rosenblatt, 1962



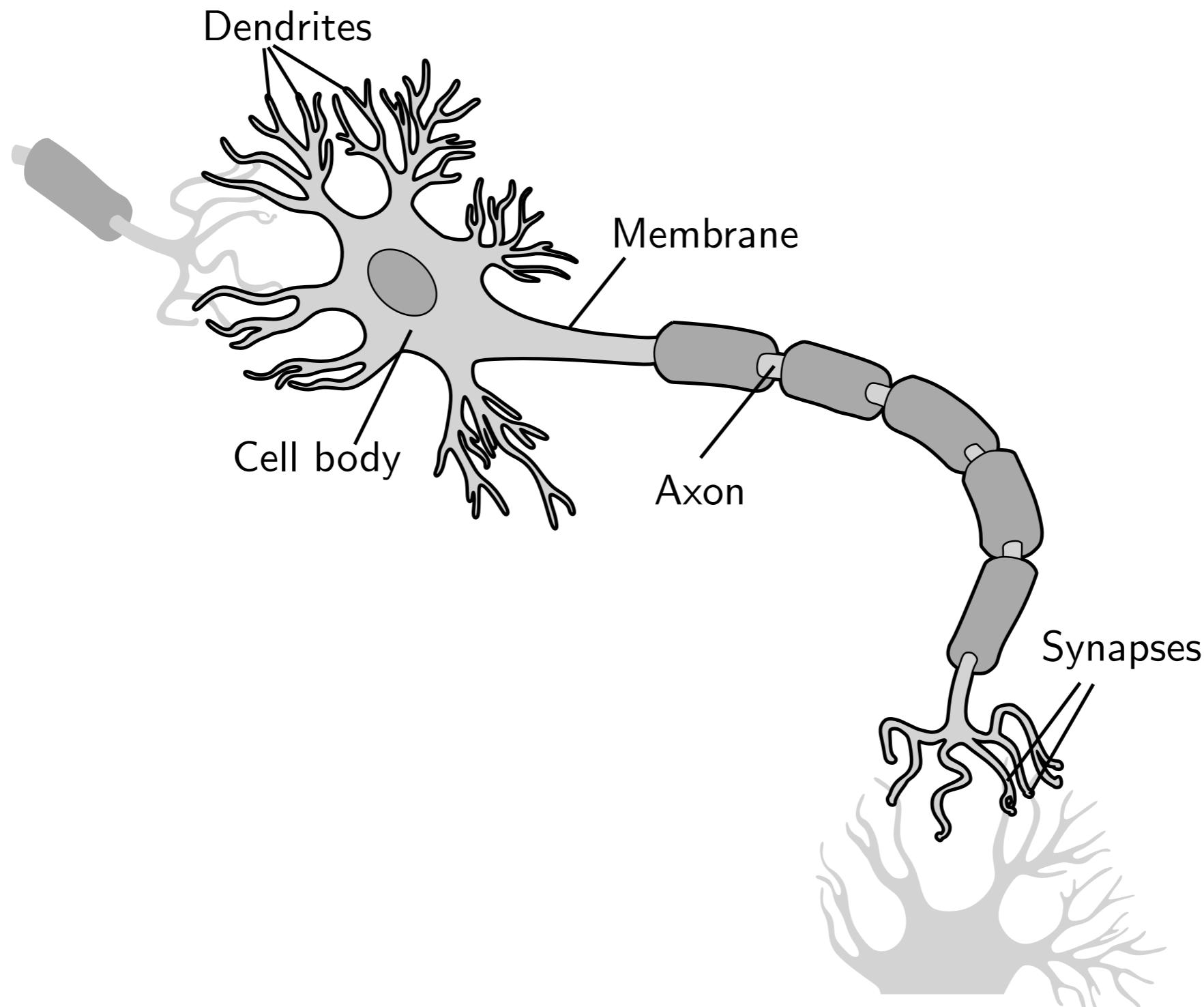
# Adaline

- Widrow & Hoff (1960s)

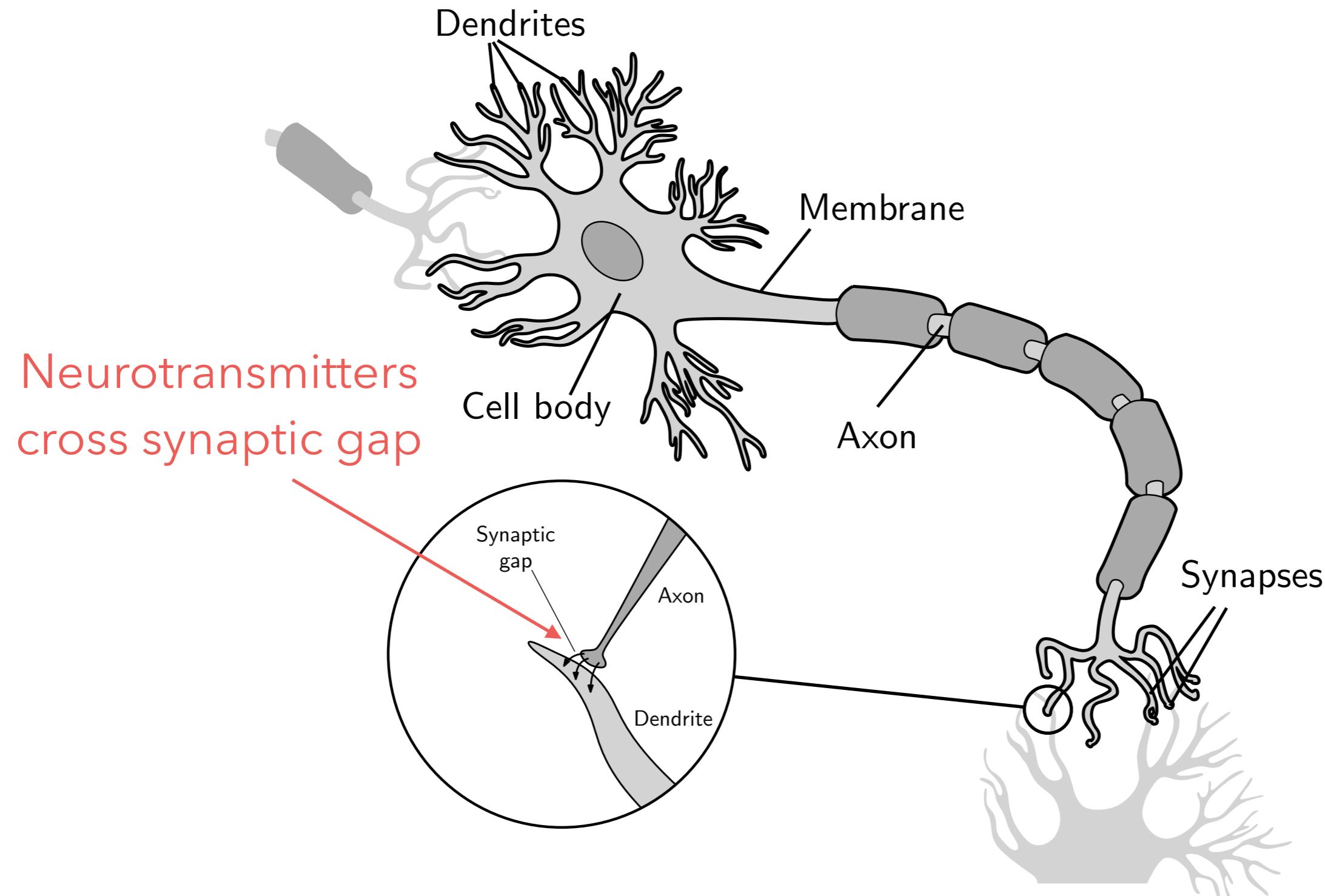


# What's this all about?

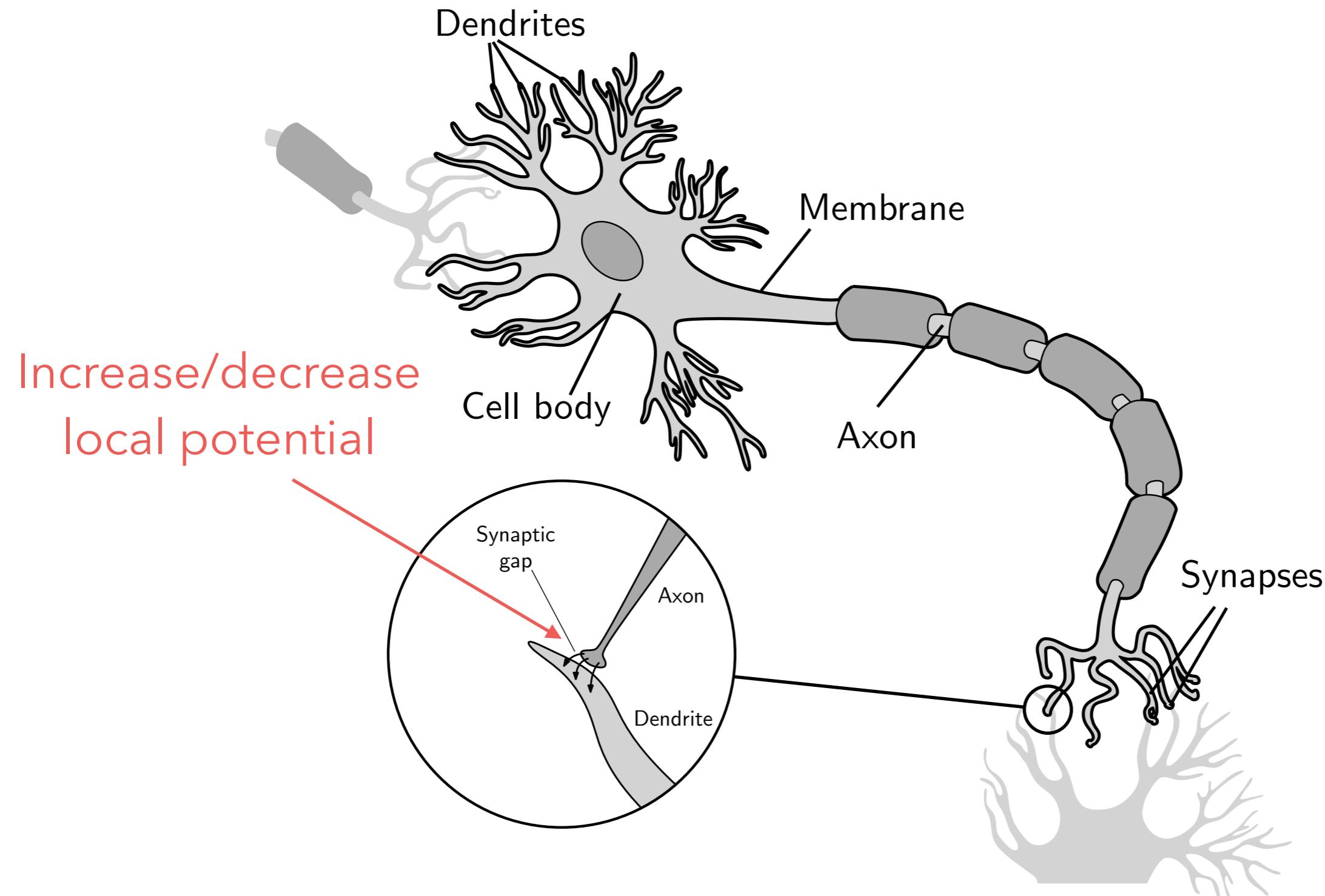
# Actual neuron



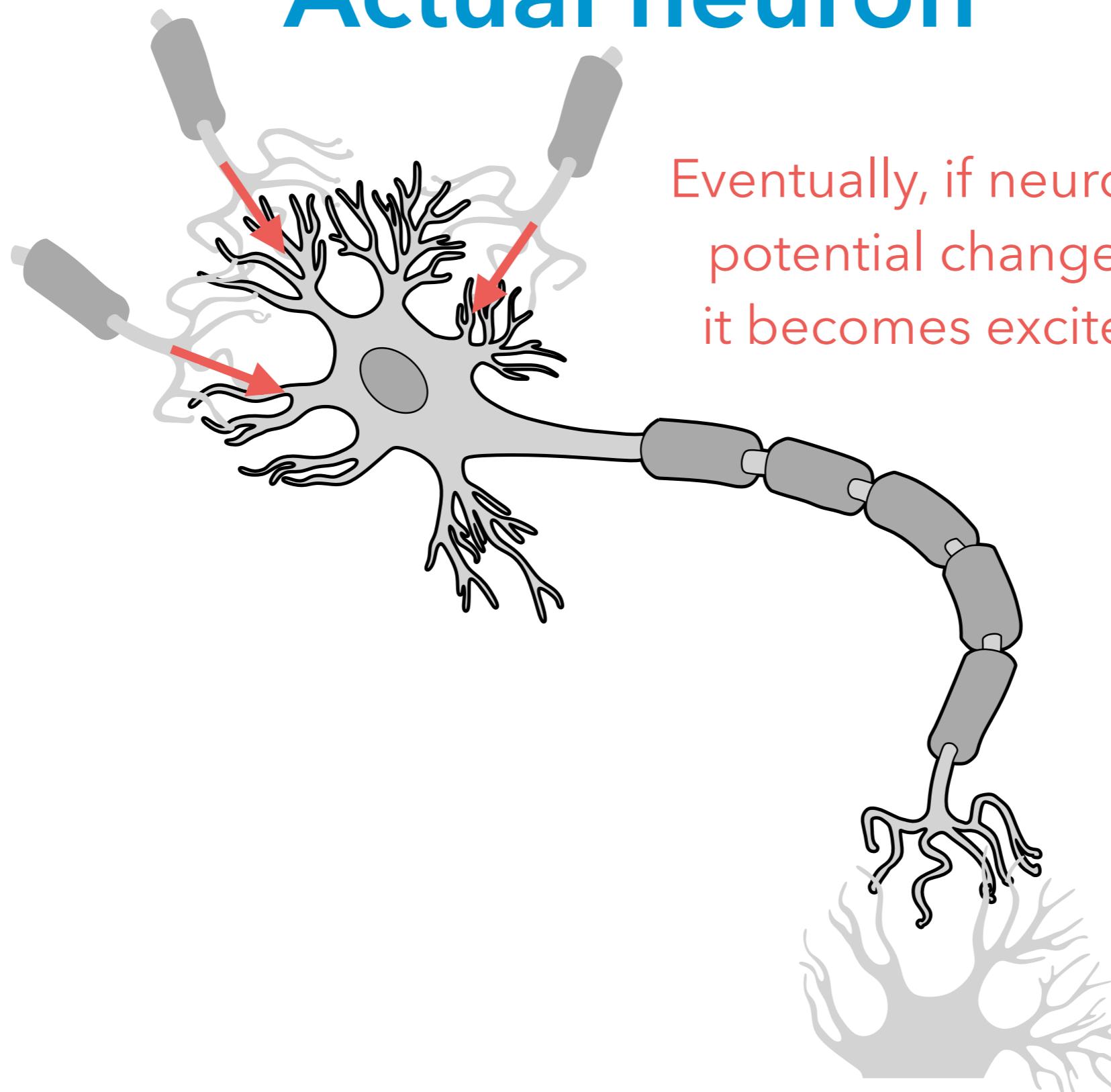
# Actual neuron



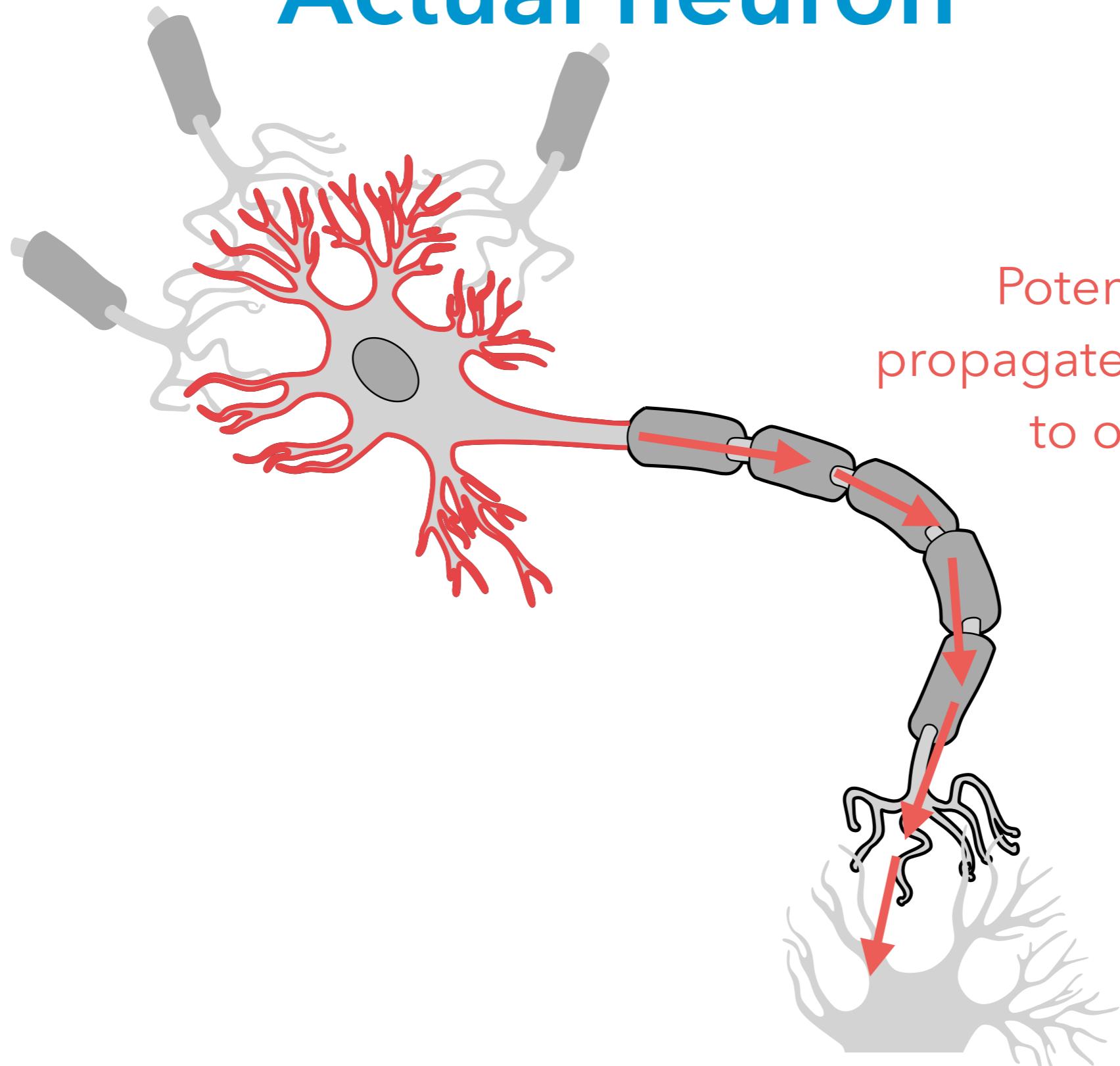
# Actual neuron



# Actual neuron

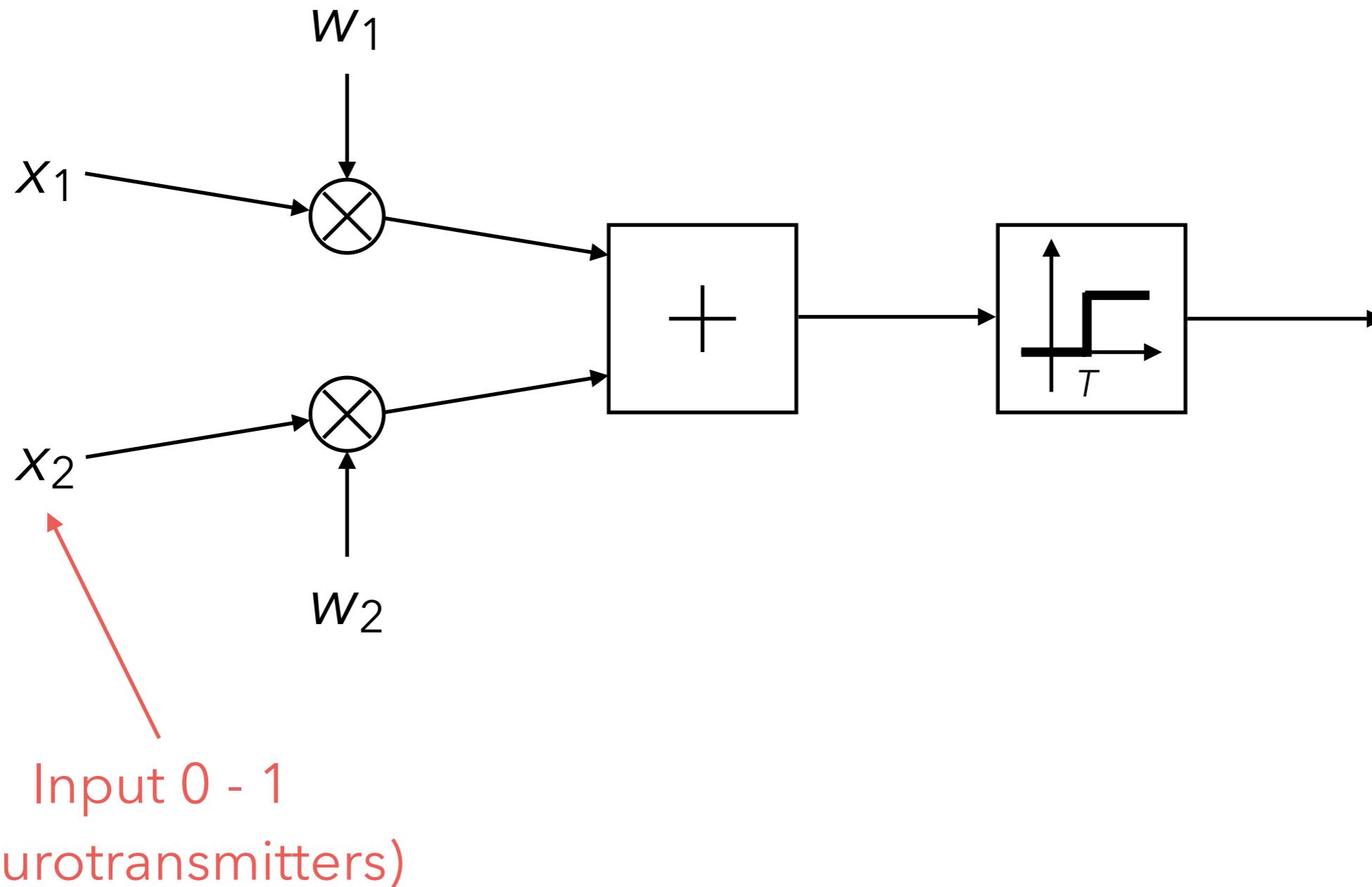


# Actual neuron

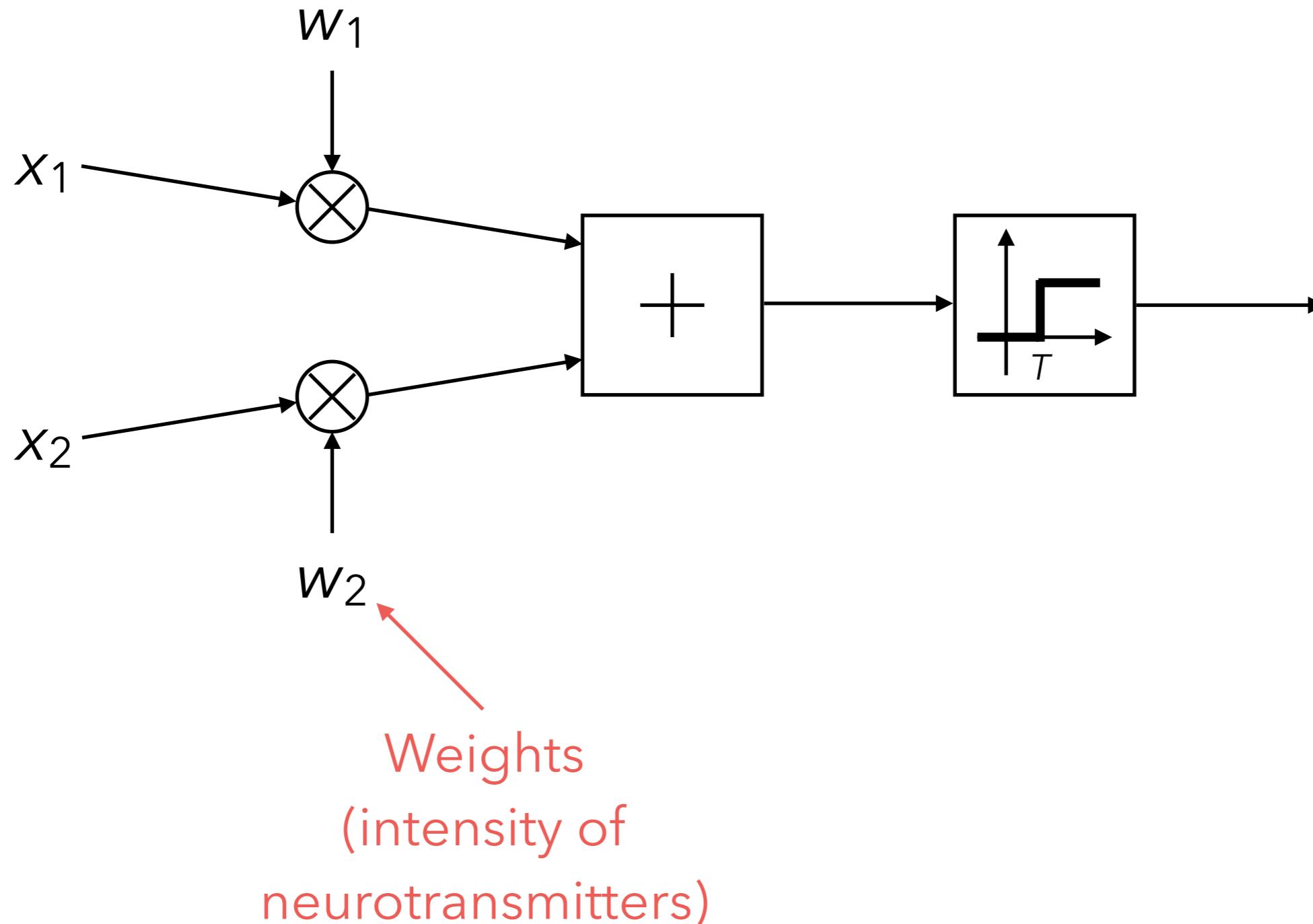


Potential is then  
propagated through axon  
to other cells

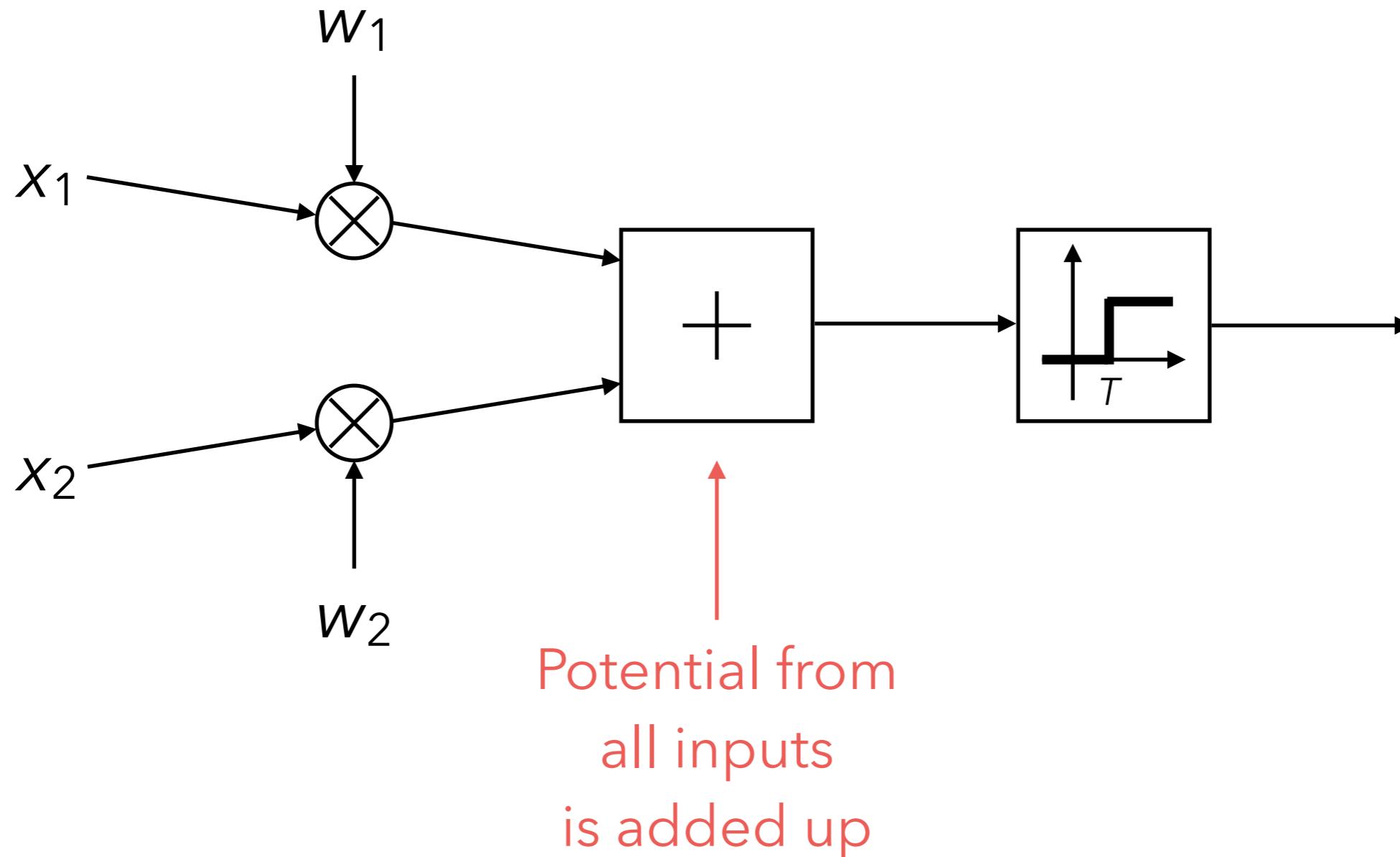
# Artificial neuron



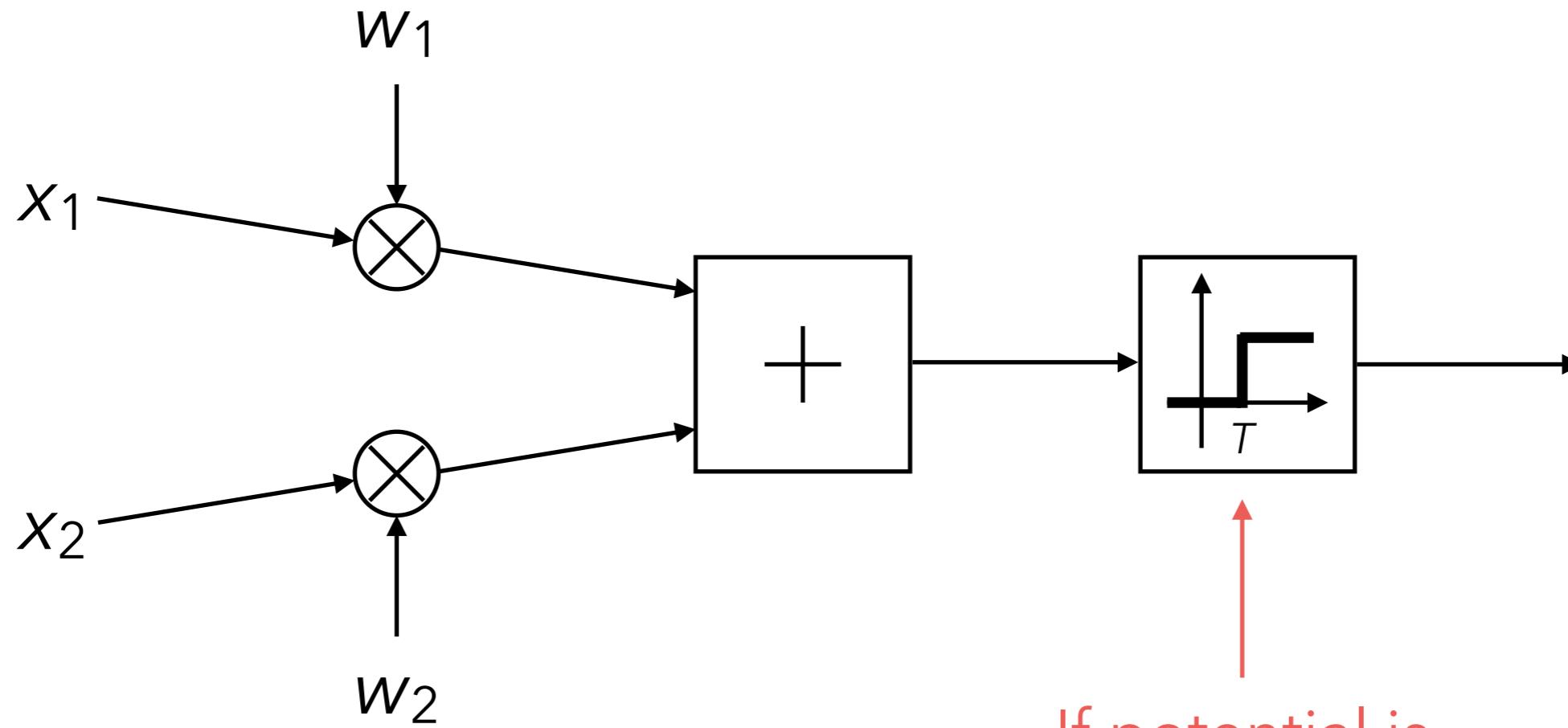
# Artificial neuron



# Artificial neuron

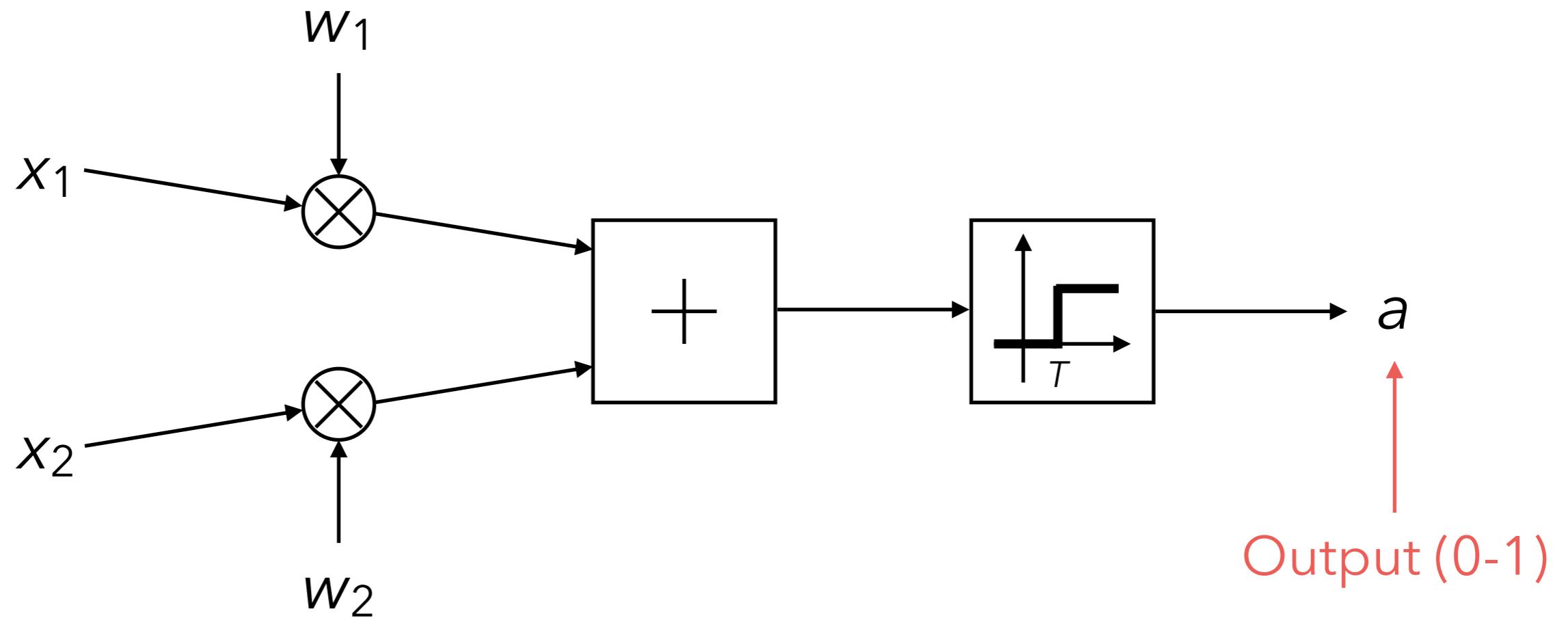


# Artificial neuron

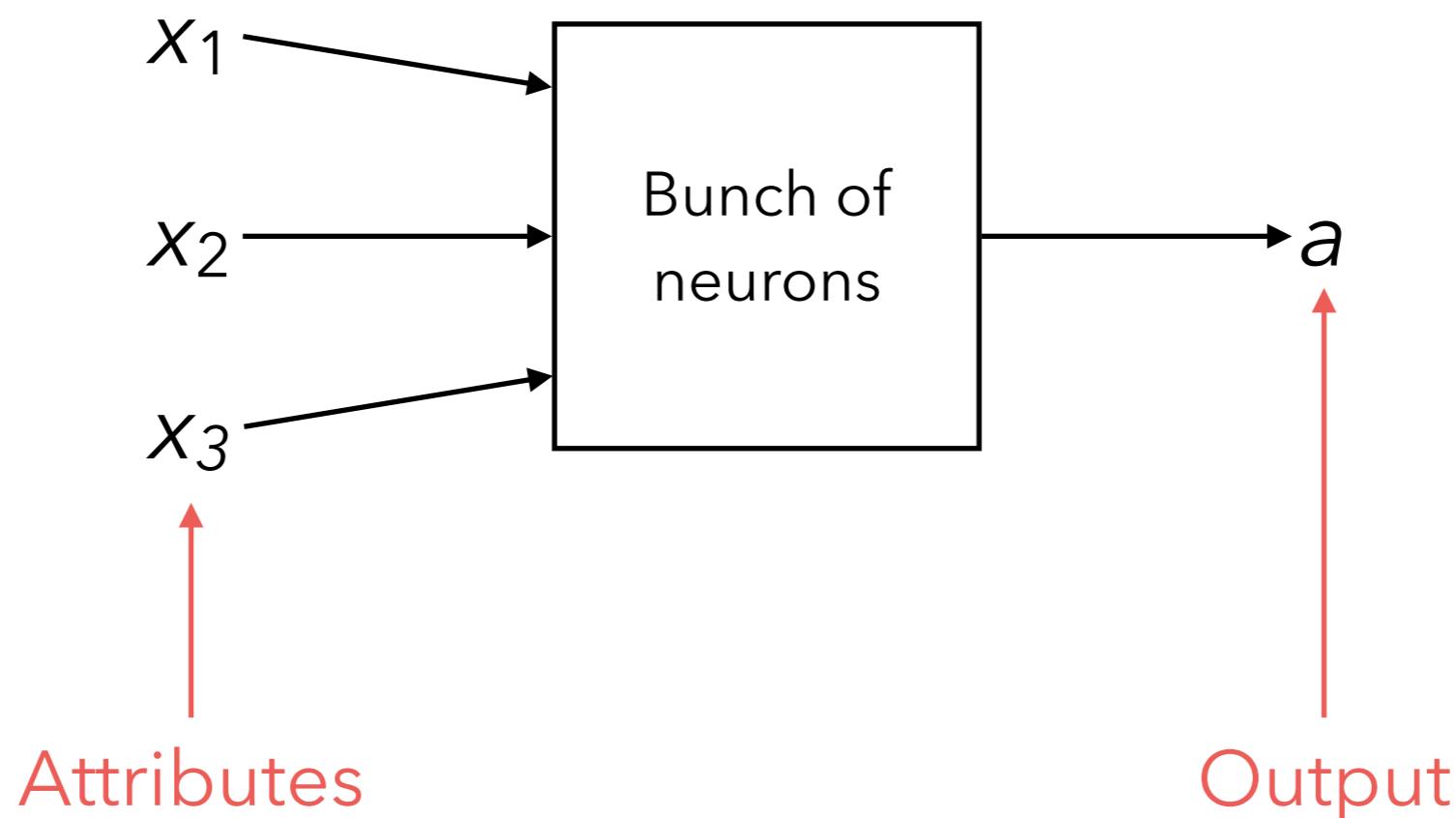


If potential is  
large enough,  
neuron “fires”

# Artificial neuron

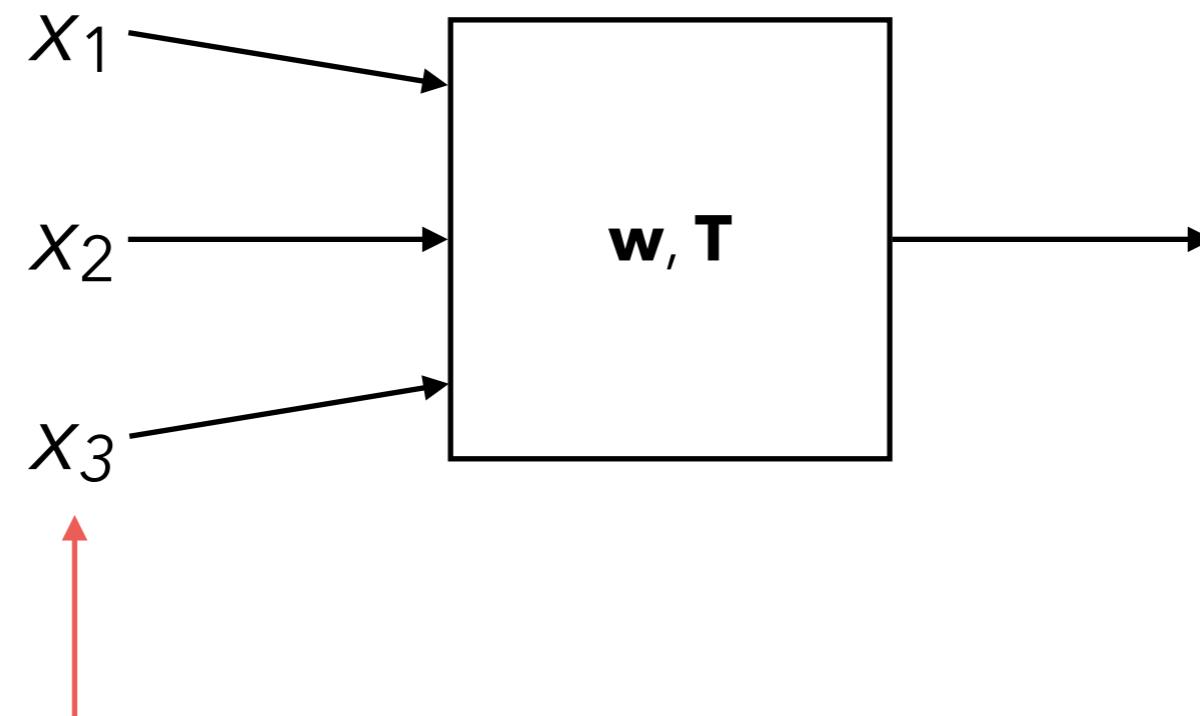


# Neural network



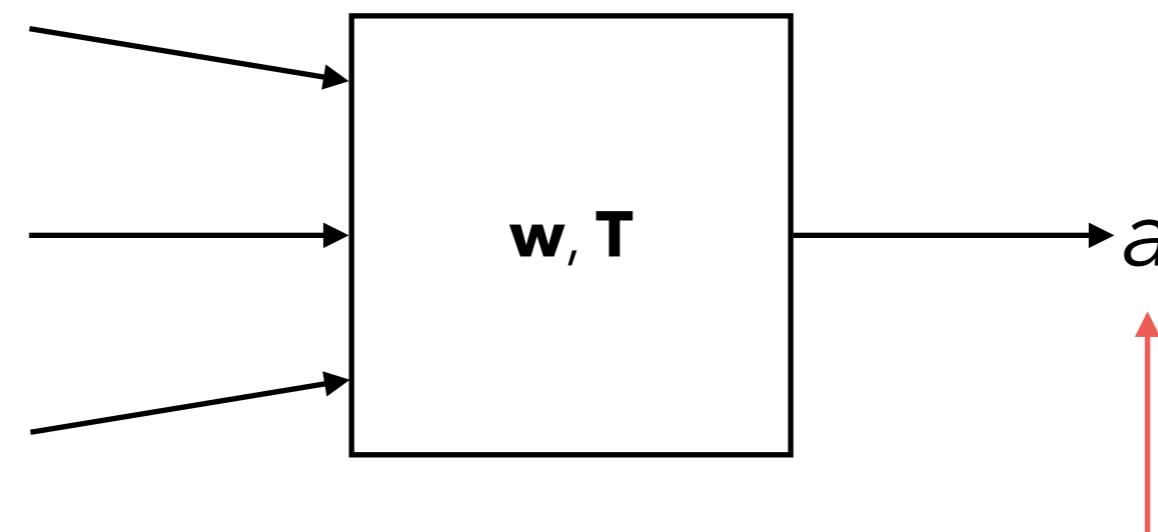
How do we learn with this?

# Neural network



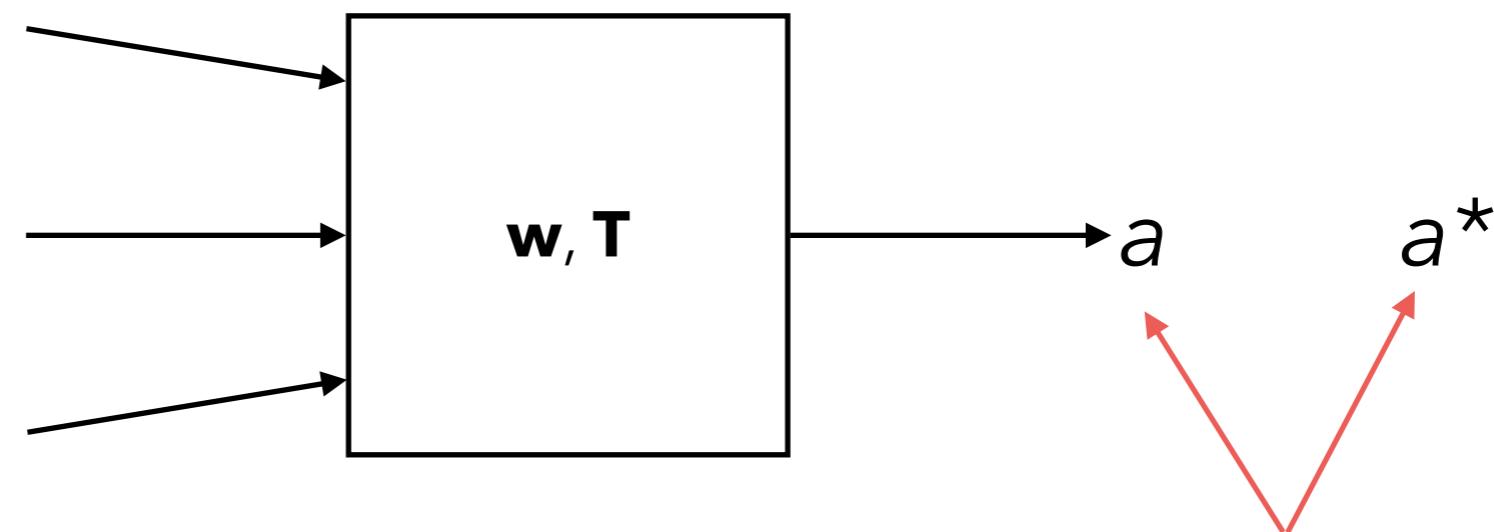
We provide  
example  $\mathbf{x}$

# Neural network



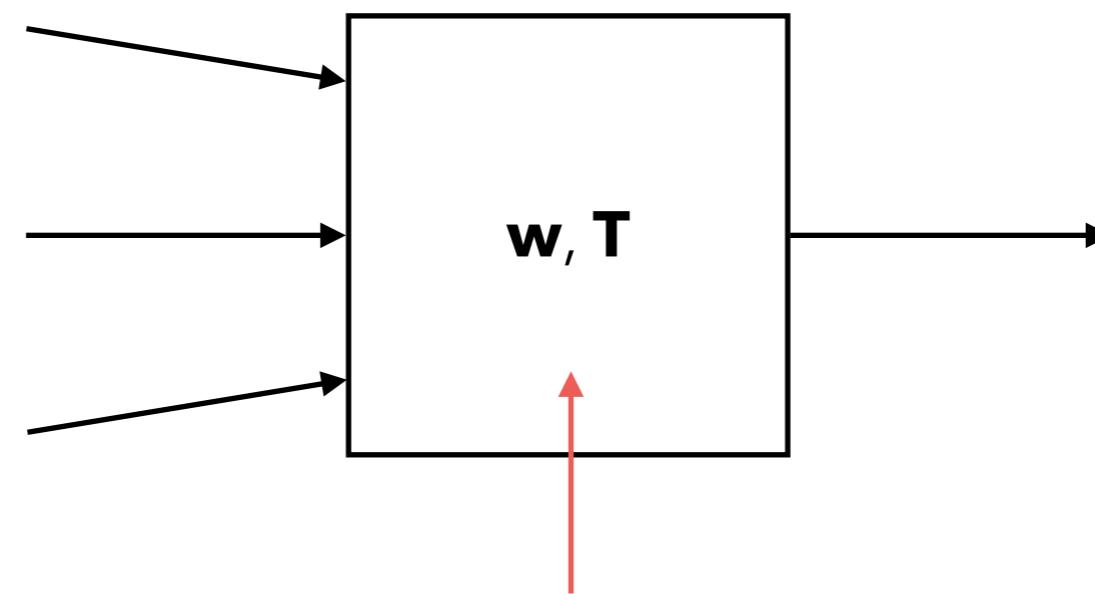
We compute the  
output,  $a = f(\mathbf{x}, \mathbf{w}, \mathbf{T})$

# Neural network



We compare the output with the desired value,  $a^*$

# Neural network



We adjust  
the weights  
and thresholds

# Training neural networks

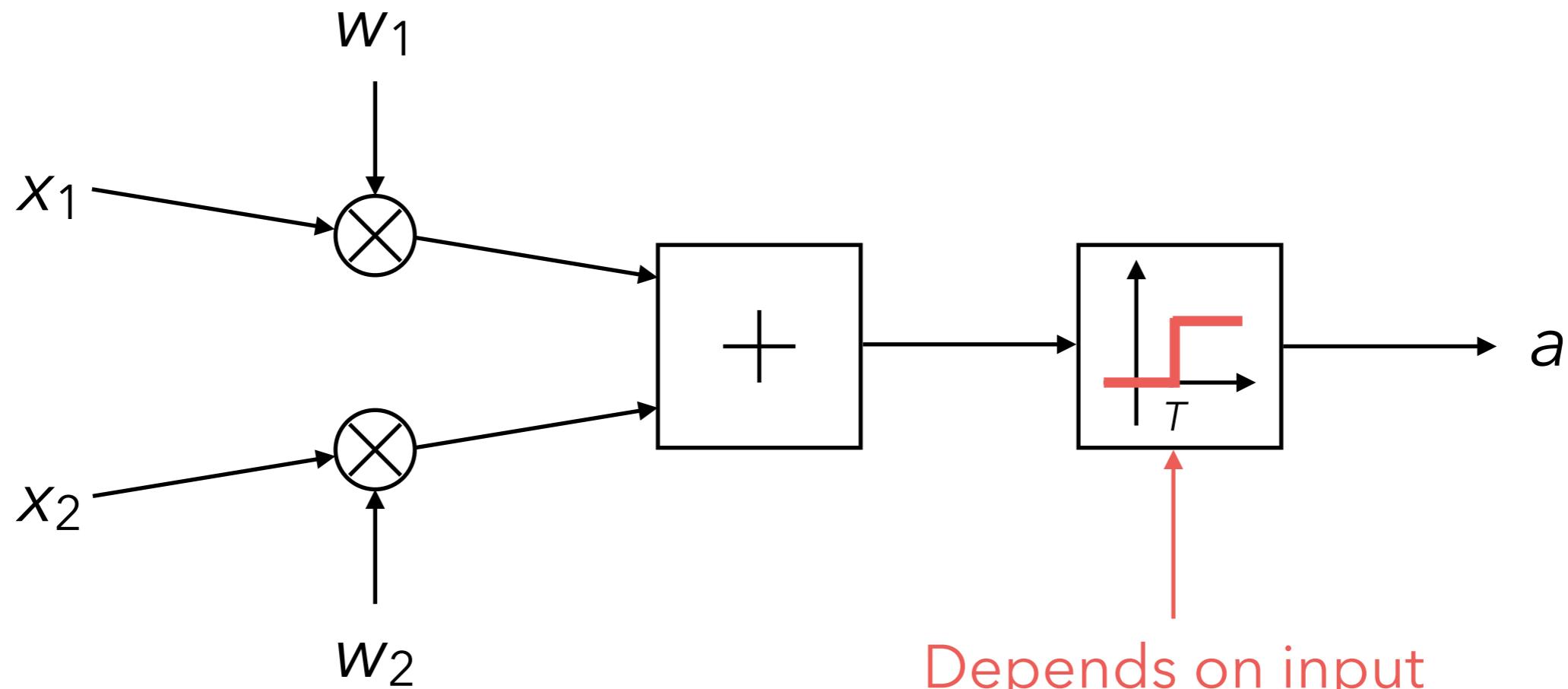
- How do we measure the performance?
  - For example, we can use the **negative log likelihood**:

$$\hat{L}_N(f) = -\frac{1}{N} \sum_{n=1}^N \log f(a_n \mid \mathbf{x}_n, \mathbf{w}, \mathbf{T})$$

- How do we adjust **w** and **T**?

... not easy. Why?

# Discontinuity!



Depends on input  
is not continuous!

$$a = \mathbb{I}(\mathbf{w}^\top \mathbf{x} > T)$$

# Step 1. Get rid of $T$

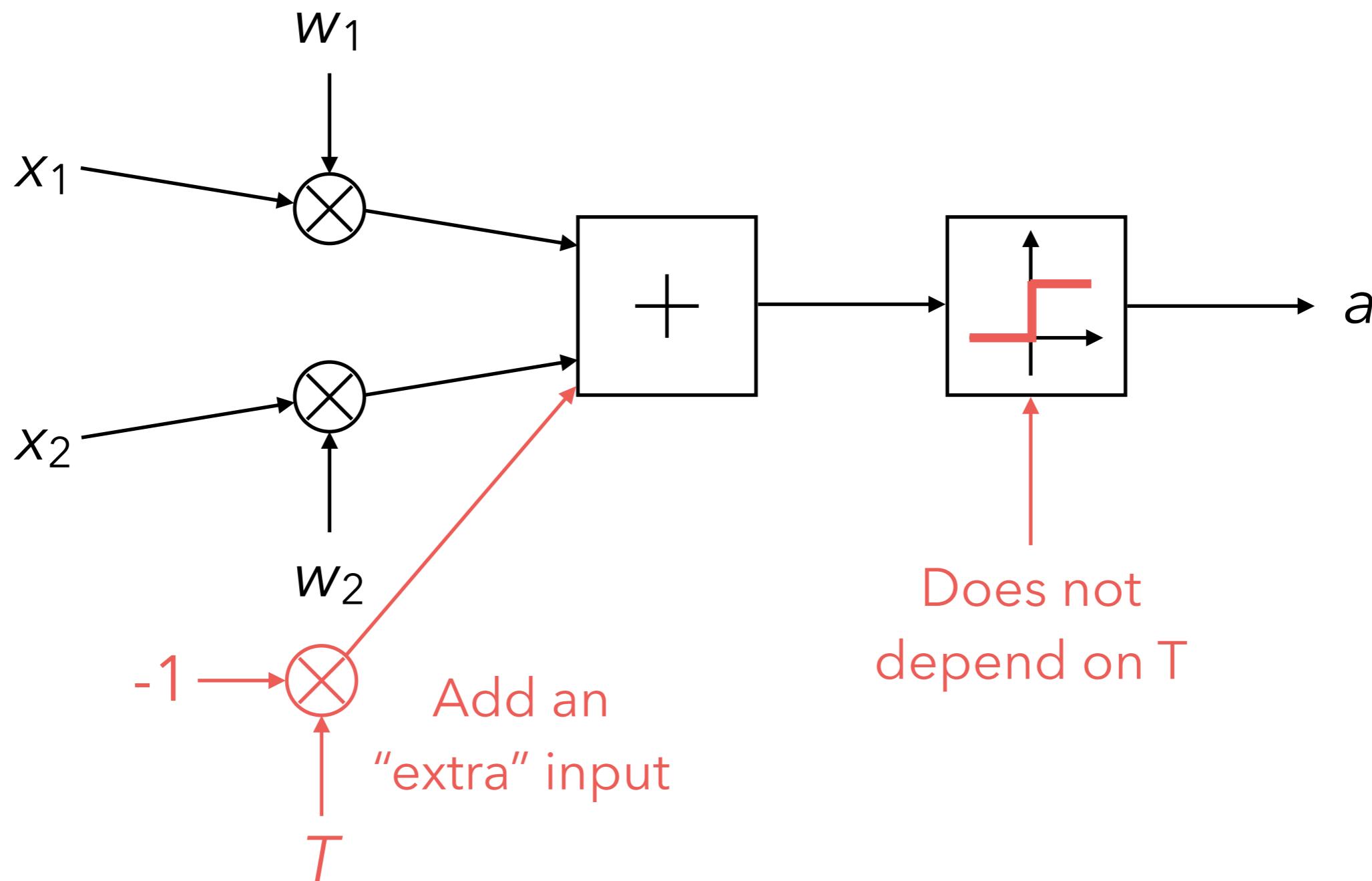
- How?
- The expression

$$a = \mathbb{I}(\mathbf{w}^\top \mathbf{x} > T)$$

is equivalent to

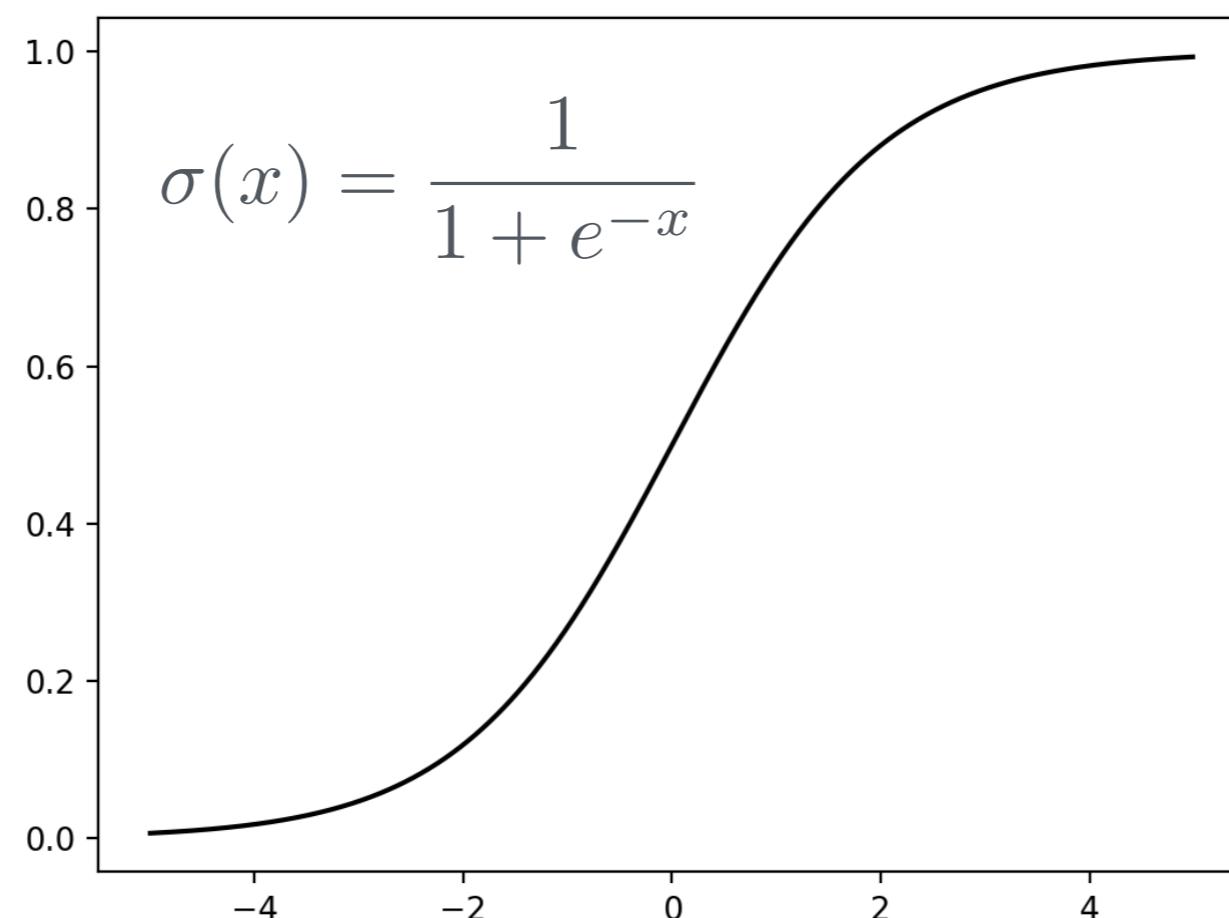
$$a = \mathbb{I}(\mathbf{w}^\top \mathbf{x} - T > 0)$$

# Step 1. Get rid of $T$



# Step 2. Smooth the threshold

- Instead of using the function  $\mathbb{I}$ , we use a smooth version of it
- For example,



# Training neural networks

- Now we can take the negative log-likelihood

$$\hat{L}_N(f) = -\frac{1}{N} \sum_{n=1}^N \log f(a_n \mid \mathbf{x}_n, \mathbf{w}, \mathbf{T})$$

and differentiate with respect to  $\mathbf{w}$  and  $\mathbf{T}$  (which is just another weight now)

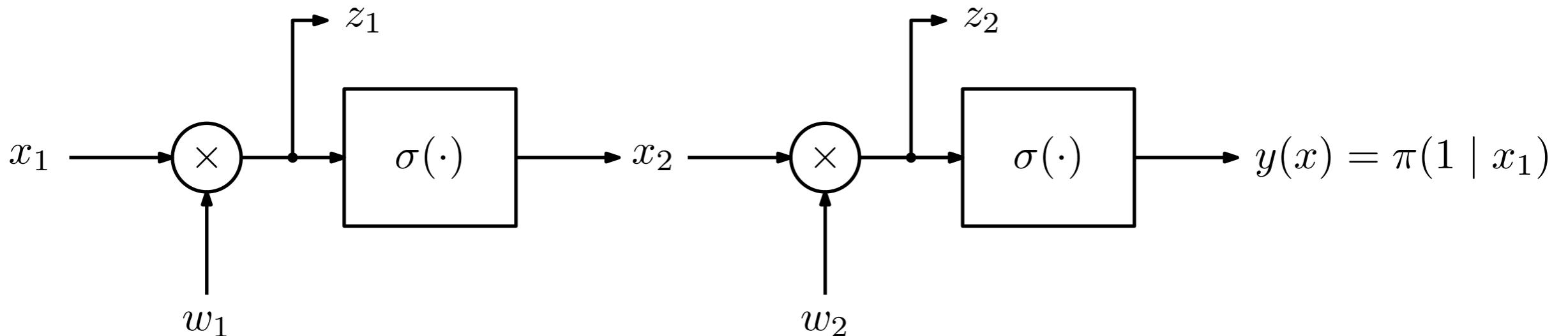
- We adjust  $\mathbf{w}$  and  $\mathbf{T}$  using **gradient descent**

$$\mathbf{w} = \mathbf{w} - \alpha \nabla_{\mathbf{w}} \hat{L}_N(f)$$



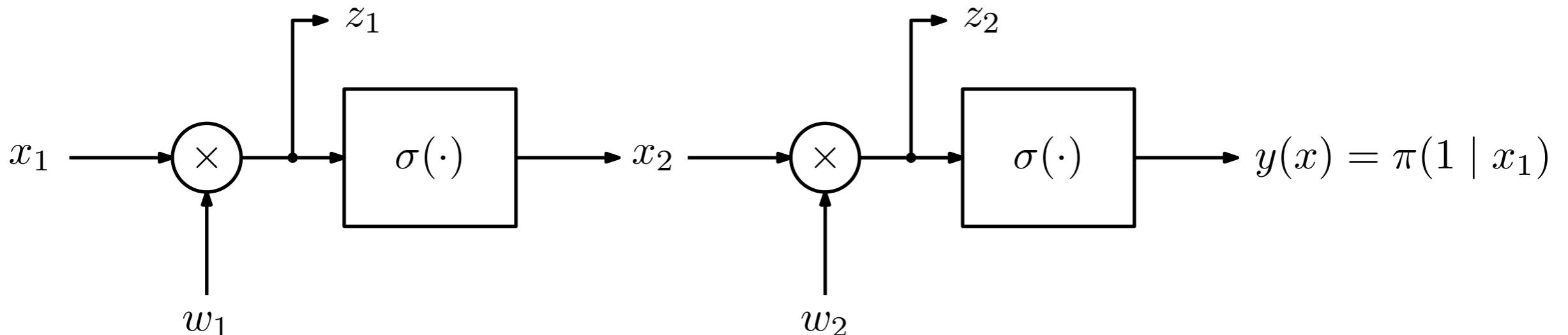
# World's simplest net

# World's simplest net



- How can we compute the dependence of the negative log-likelihood on  $w_1$  and  $w_2$ ?
  - We use the chain rule (of derivatives)

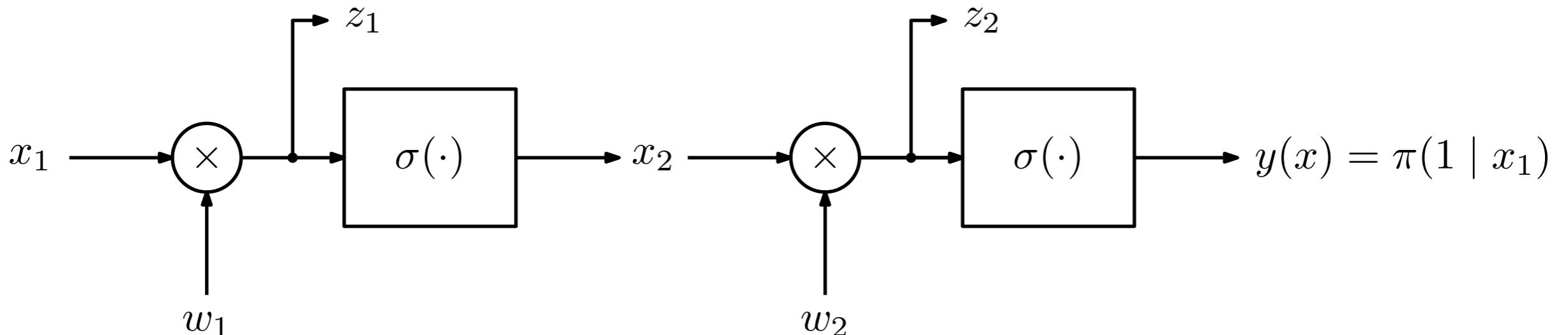
# World's simplest net



$$\frac{\partial \hat{L}_N}{\partial w_1} = \frac{\partial \hat{L}_N}{\partial y} \times \frac{\partial y}{\partial w_1}$$

$$\frac{\partial \hat{L}_N}{\partial w_2} = \frac{\partial \hat{L}_N}{\partial y} \times \frac{\partial y}{\partial w_2}$$

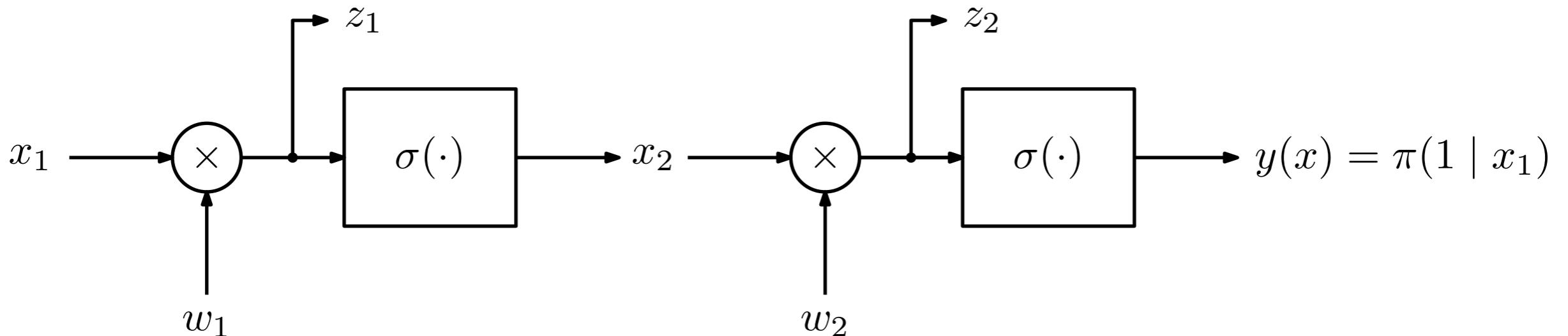
# World's simplest net



$$\frac{\partial \hat{L}_N}{\partial w_1} = \frac{\partial \hat{L}_N}{\partial y} \times \frac{\partial y}{\partial z_2} \times \frac{\partial z_2}{\partial w_1}$$

$$\frac{\partial \hat{L}_N}{\partial w_2} = \frac{\partial \hat{L}_N}{\partial y} \times \frac{\partial y}{\partial z_2} \times \frac{\partial z_2}{\partial w_2}$$
✓

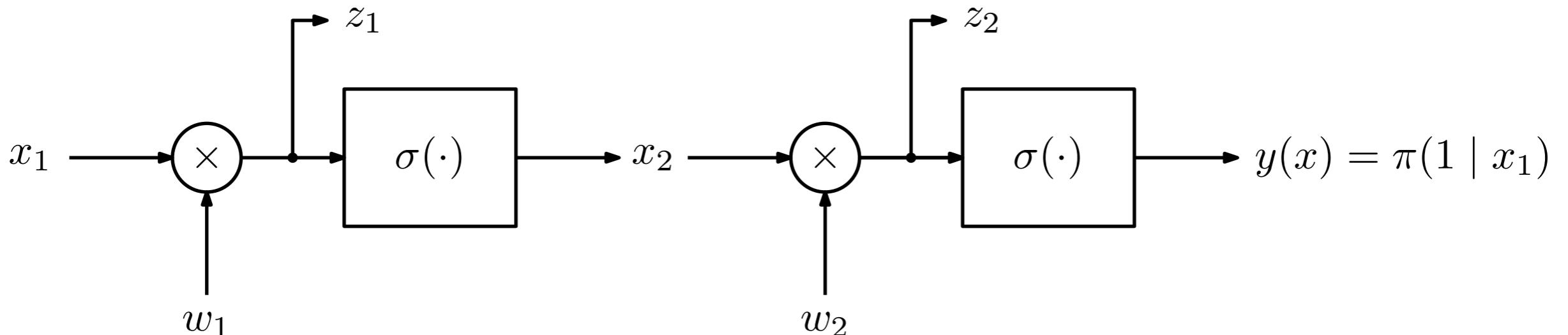
# World's simplest net



$$\frac{\partial \hat{L}_N}{\partial w_1} = \frac{\partial \hat{L}_N}{\partial y} \times \frac{\partial y}{\partial z_2} \times \frac{\partial z_2}{\partial x_2} \times \frac{\partial x_2}{\partial w_1}$$

$$\frac{\partial \hat{L}_N}{\partial w_2} = \frac{\partial \hat{L}_N}{\partial y} \times \frac{\partial y}{\partial z_2} \times \frac{\partial z_2}{\partial w_2}$$

# World's simplest net

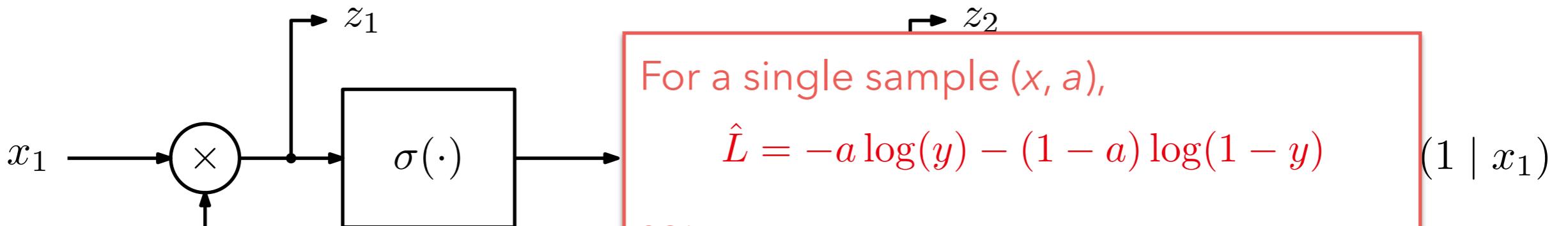


$$\frac{\partial \hat{L}_N}{\partial w_1} = \frac{\partial \hat{L}_N}{\partial y} \times \frac{\partial y}{\partial z_2} \times \frac{\partial z_2}{\partial x_2} \times \frac{\partial x_2}{\partial z_1} \times \frac{\partial z_1}{\partial w_1} \quad \checkmark$$

$$\frac{\partial \hat{L}_N}{\partial w_2} = \frac{\partial \hat{L}_N}{\partial y} \times \frac{\partial y}{\partial z_2} \times \frac{\partial z_2}{\partial w_2}$$

Now let's break  
these down!

# World's simplest net



For a single sample  $(x, a)$ ,

$$\hat{L} = -a \log(y) - (1 - a) \log(1 - y) \quad (1 | x_1)$$

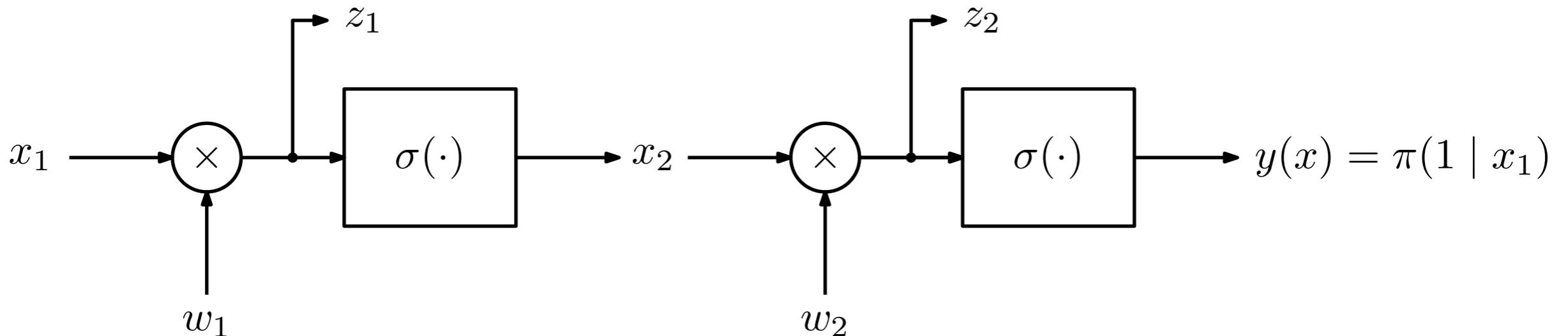
so:

$$\frac{\partial \hat{L}}{\partial y} = \frac{-a}{y} + \frac{1 - a}{1 - y}$$

$$\frac{\partial \hat{L}_N}{\partial w_1} = \boxed{\frac{\partial \hat{L}_N}{\partial y}} \times \frac{\partial y}{\partial z_2} \times \frac{\partial z_2}{\partial x_2} \times \frac{\partial x_2}{\partial z_1} \times \frac{\partial z_1}{\partial w_1}$$

$$\frac{\partial \hat{L}_N}{\partial w_2} = \boxed{\frac{\partial \hat{L}_N}{\partial y}} \times \frac{\partial y}{\partial z_2} \times \frac{\partial z_2}{\partial w_2}$$

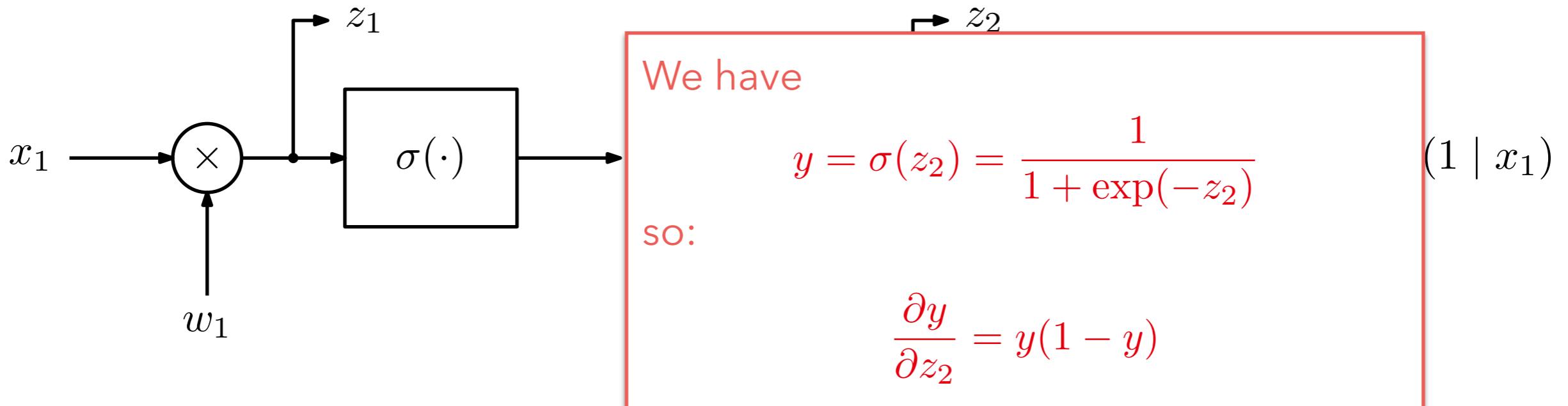
# World's simplest net



$$\frac{\partial \hat{L}_N}{\partial w_1} = \left( \frac{-a}{y} + \frac{1-a}{1-y} \right) \times \frac{\partial y}{\partial z_2} \times \frac{\partial z_2}{\partial x_2} \times \frac{\partial x_2}{\partial z_1} \times \frac{\partial z_1}{\partial w_1}$$

$$\frac{\partial \hat{L}_N}{\partial w_2} = \left( \frac{-a}{y} + \frac{1-a}{1-y} \right) \times \frac{\partial y}{\partial z_2} \times \frac{\partial z_2}{\partial w_2}$$

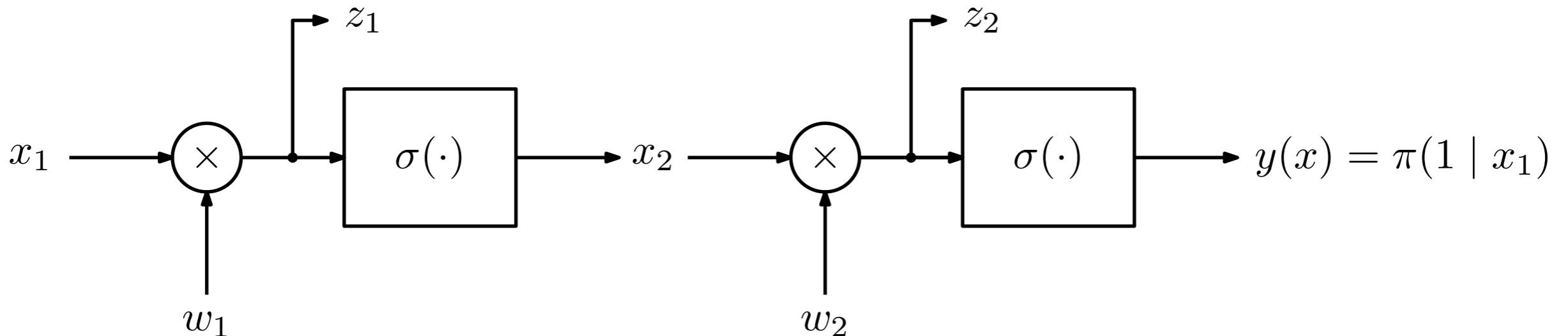
# World's simplest net



$$\frac{\partial \hat{L}_N}{\partial w_1} = \left( \frac{-a}{y} + \frac{1-a}{1-y} \right) \times \boxed{\frac{\partial y}{\partial z_2}} \times \frac{\partial z_2}{\partial x_2} \times \frac{\partial x_2}{\partial z_1} \times \frac{\partial z_1}{\partial w_1}$$

$$\frac{\partial \hat{L}_N}{\partial w_2} = \left( \frac{-a}{y} + \frac{1-a}{1-y} \right) \times \boxed{\frac{\partial y}{\partial z_2}} \times \frac{\partial z_2}{\partial w_2}$$

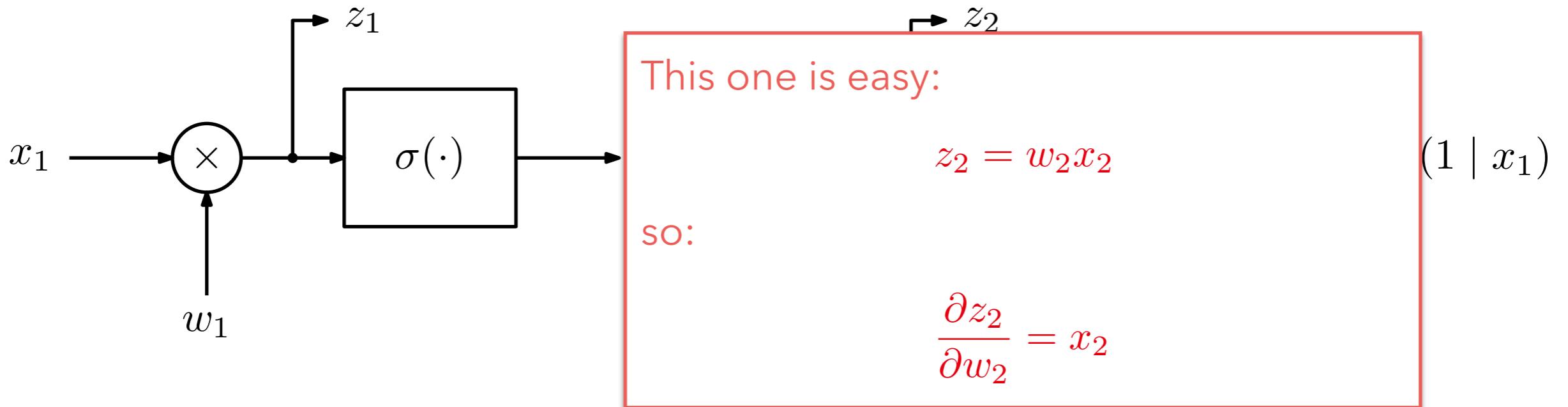
# World's simplest net



$$\frac{\partial \hat{L}_N}{\partial w_1} = (y - a) \times \frac{\partial z_2}{\partial x_2} \times \frac{\partial x_2}{\partial z_1} \times \frac{\partial z_1}{\partial w_1}$$

$$\frac{\partial \hat{L}_N}{\partial w_2} = (y - a) \times \frac{\partial z_2}{\partial w_2}$$

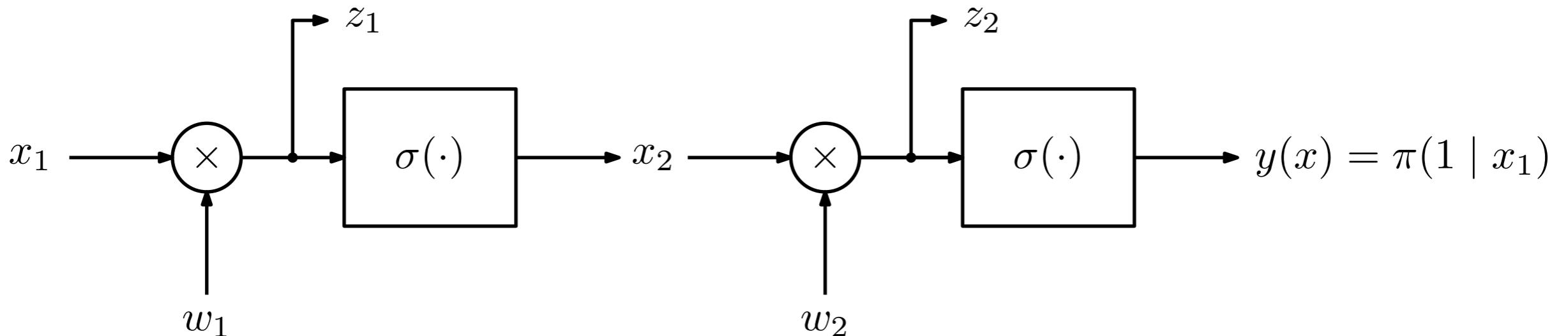
# World's simplest net



$$\frac{\partial \hat{L}_N}{\partial w_1} = (y - a) \times \frac{\partial z_2}{\partial x_2} \times \frac{\partial x_2}{\partial z_1} \times \frac{\partial z_1}{\partial w_1}$$

$$\frac{\partial \hat{L}_N}{\partial w_2} = (y - a) \times \boxed{\frac{\partial z_2}{\partial w_2}}$$

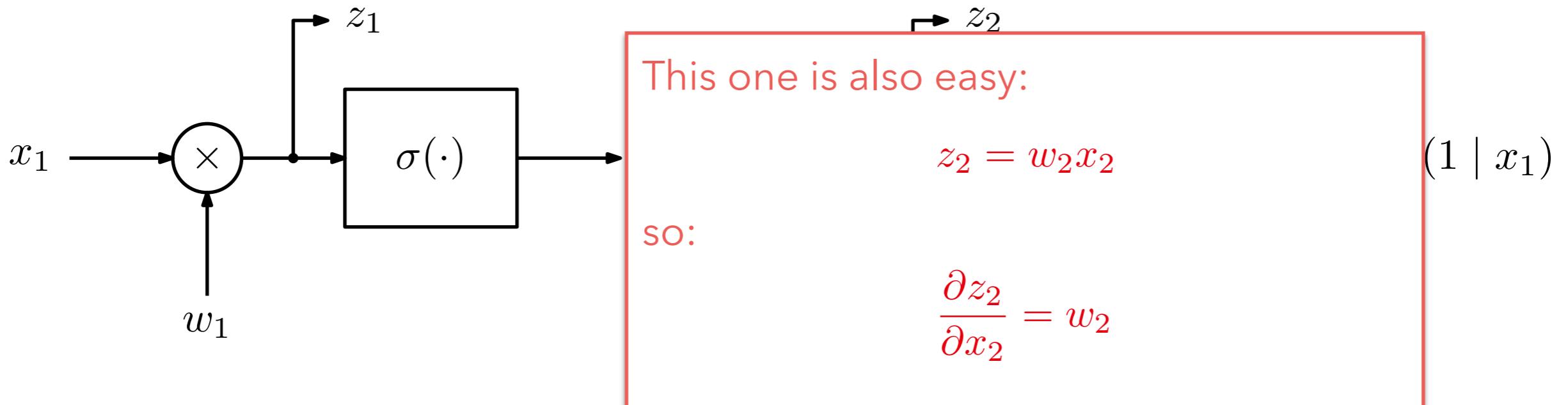
# World's simplest net



$$\frac{\partial \hat{L}_N}{\partial w_1} = (y - a) \times \frac{\partial z_2}{\partial x_2} \times \frac{\partial x_2}{\partial z_1} \times \frac{\partial z_1}{\partial w_1}$$

$$\frac{\partial \hat{L}_N}{\partial w_2} = (y - a)x_2$$

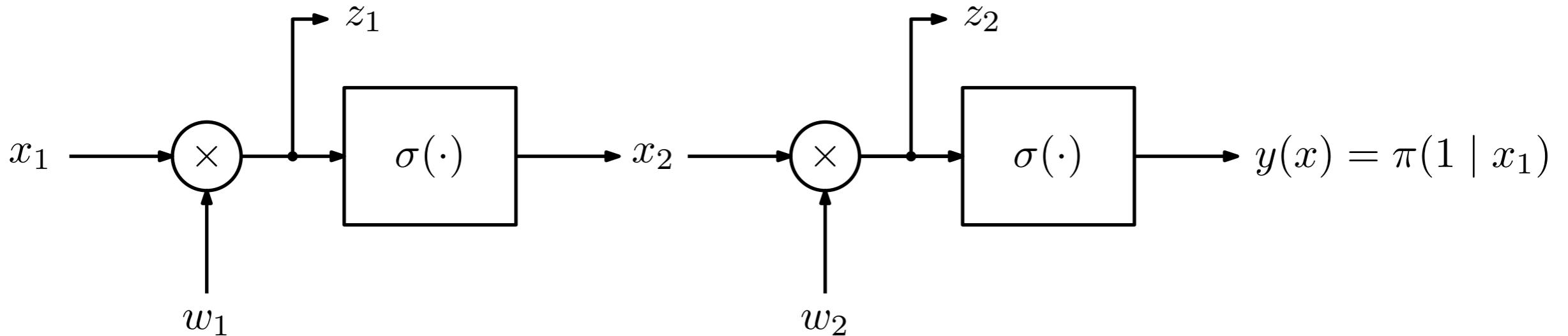
# World's simplest net



$$\frac{\partial \hat{L}_N}{\partial w_1} = (y - a) \times \boxed{\frac{\partial z_2}{\partial x_2}} \times \frac{\partial x_2}{\partial z_1} \times \frac{\partial z_1}{\partial w_1}$$

$$\frac{\partial \hat{L}_N}{\partial w_2} = (y - a)x_2$$

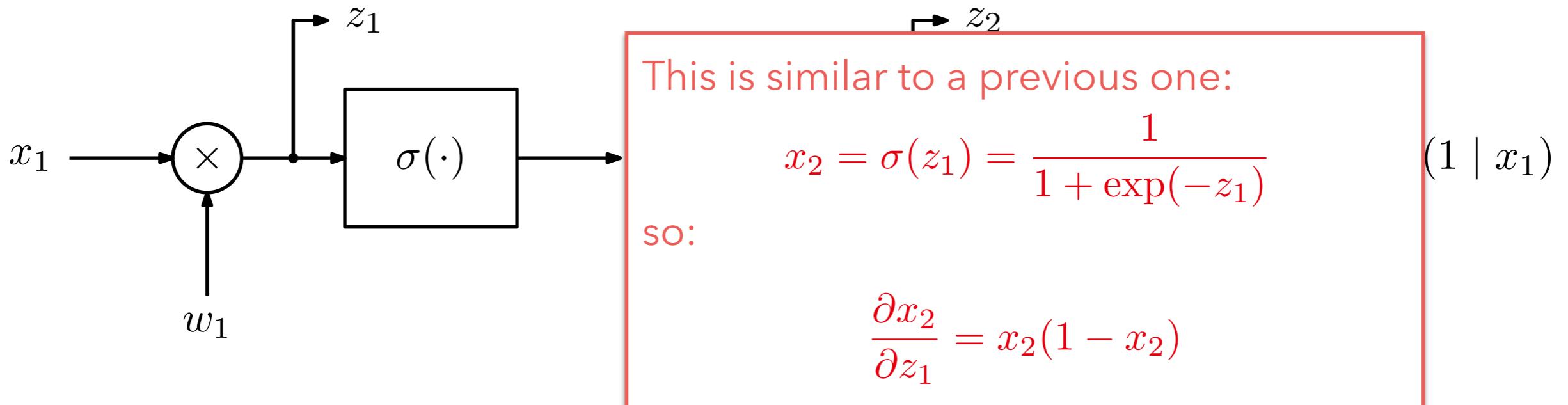
# World's simplest net



$$\frac{\partial \hat{L}_N}{\partial w_1} = (y - a)w_2 \times \frac{\partial x_2}{\partial z_1} \times \frac{\partial z_1}{\partial w_1}$$

$$\frac{\partial \hat{L}_N}{\partial w_2} = (y - a)x_2$$

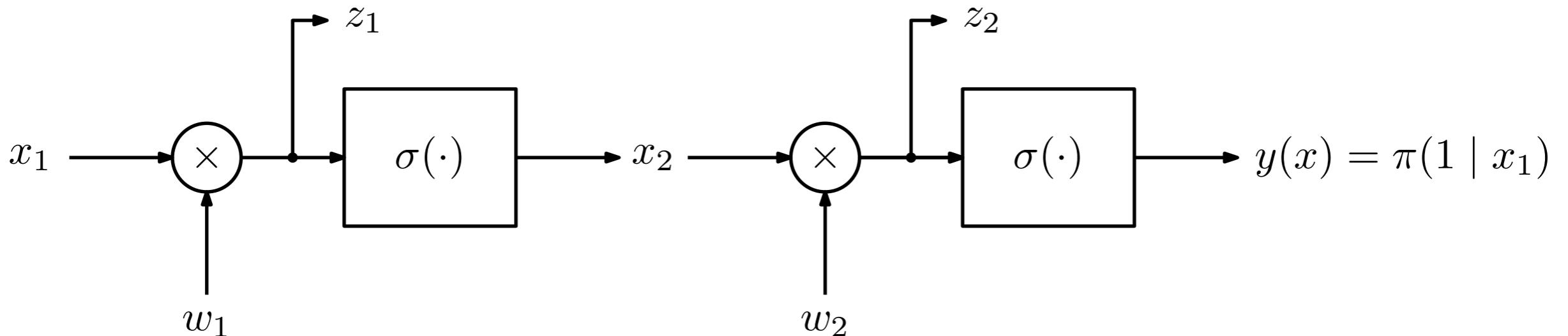
# World's simplest net



$$\frac{\partial \hat{L}_N}{\partial w_1} = (y - a)w_2 \times \boxed{\frac{\partial x_2}{\partial z_1}} \times \frac{\partial z_1}{\partial w_1}$$

$$\frac{\partial \hat{L}_N}{\partial w_2} = (y - a)x_2$$

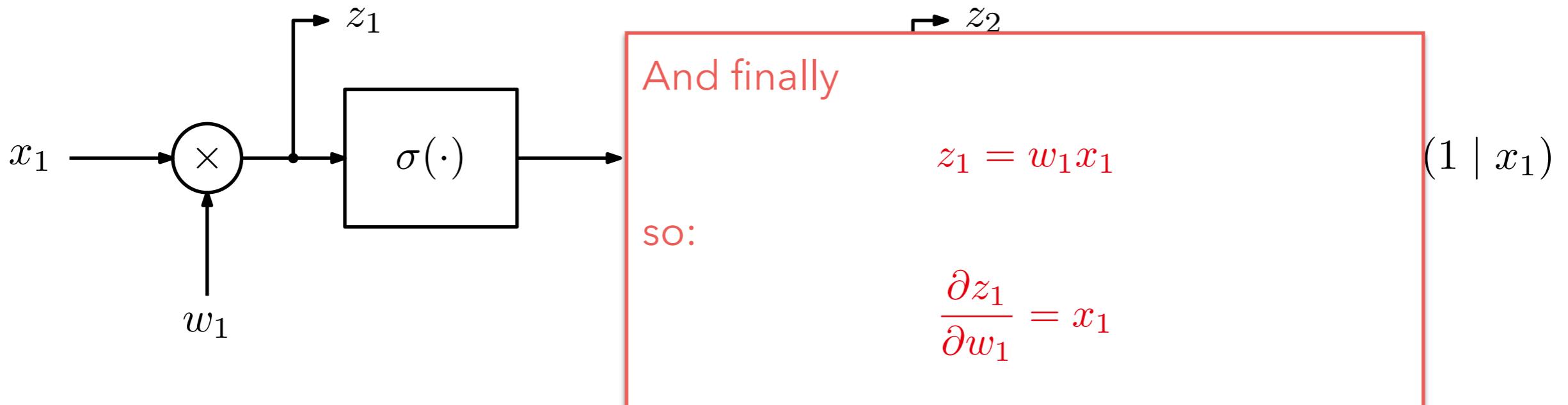
# World's simplest net



$$\frac{\partial \hat{L}_N}{\partial w_1} = (y - a)w_2x_2(1 - x_2) \times \frac{\partial z_1}{\partial w_1}$$

$$\frac{\partial \hat{L}_N}{\partial w_2} = (y - a)x_2$$

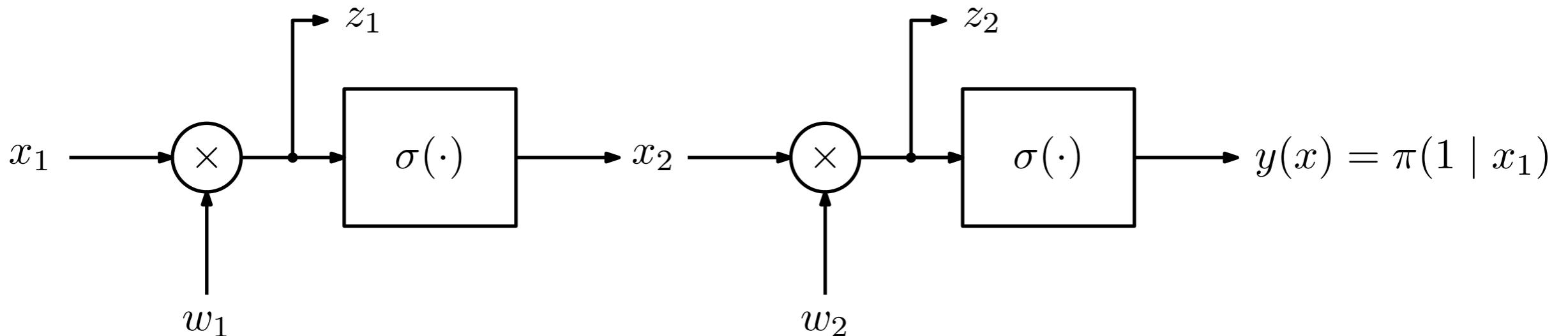
# World's simplest net



$$\frac{\partial \hat{L}_N}{\partial w_1} = (y - a)w_2x_2(1 - x_2) \times \boxed{\frac{\partial z_1}{\partial w_1}}$$

$$\frac{\partial \hat{L}_N}{\partial w_2} = (y - a)x_2$$

# World's simplest net

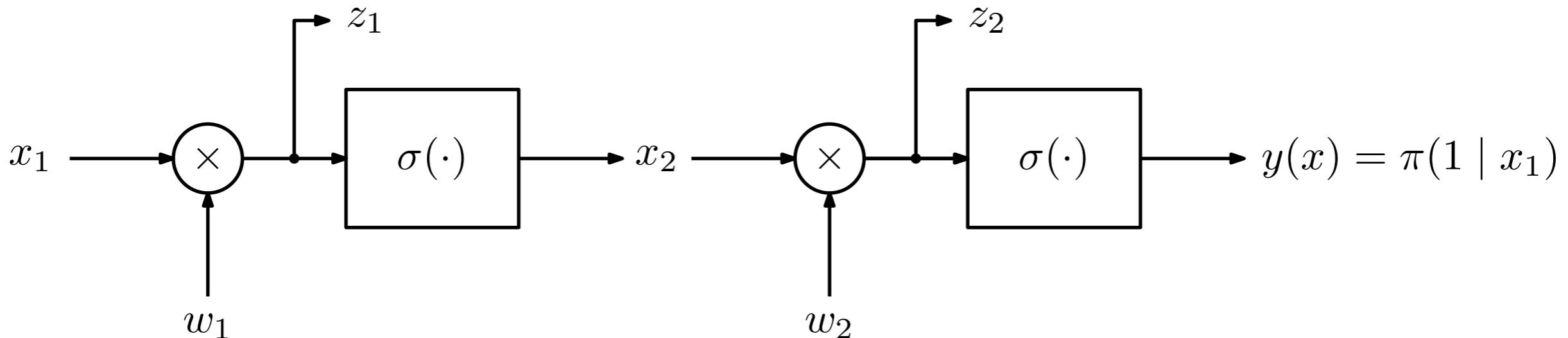


$$\frac{\partial \hat{L}_N}{\partial w_1} = (y - a)w_2x_2(1 - x_2)x_1$$

$$\frac{\partial \hat{L}_N}{\partial w_2} = (y - a)x_2$$

Looking more carefully...

# World's simplest net

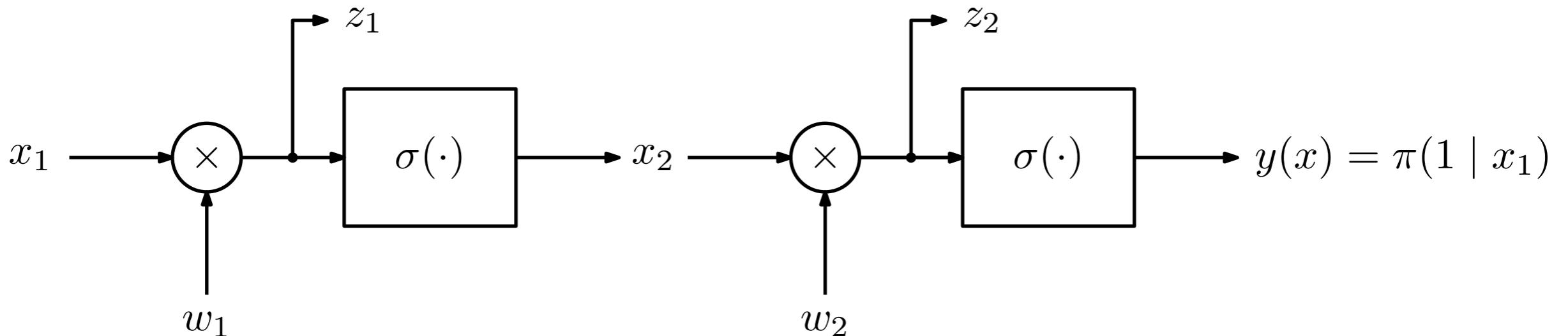


$$\frac{\partial \hat{L}_N}{\partial w_1} = (y - a)w_2x_2(1 - x_2)x_1$$

$$\frac{\partial \hat{L}_N}{\partial w_2} = (y - a)x_2$$

The derivative with respect to  $w_i$  is something  $\times x_i$

# World's simplest net

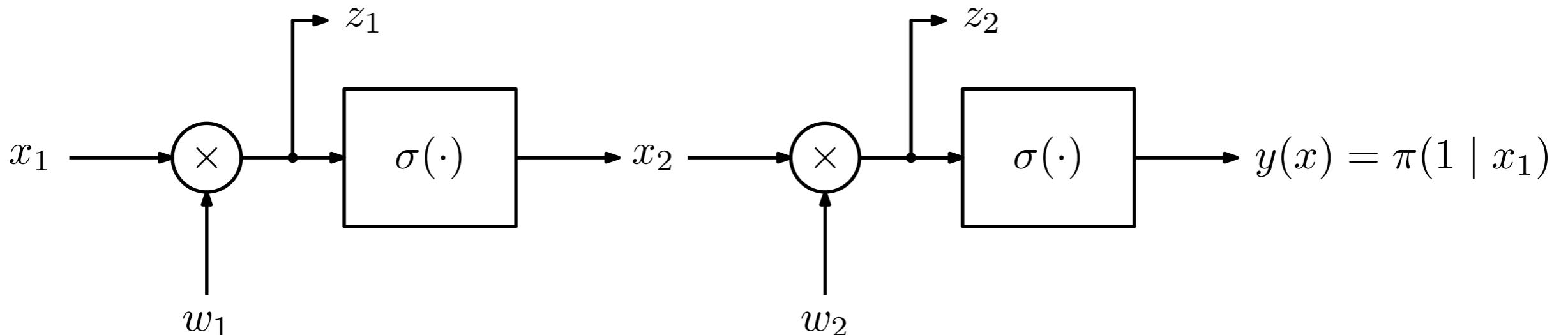


$$\frac{\partial \hat{L}_N}{\partial w_1} = \frac{(y - a)w_2x_2(1 - x_2)}{\delta_1}x_1$$

$$\frac{\partial \hat{L}_N}{\partial w_2} = \frac{(y - a)x_2}{\delta_2}$$

We will call that something  $\delta_i$

# World's simplest net



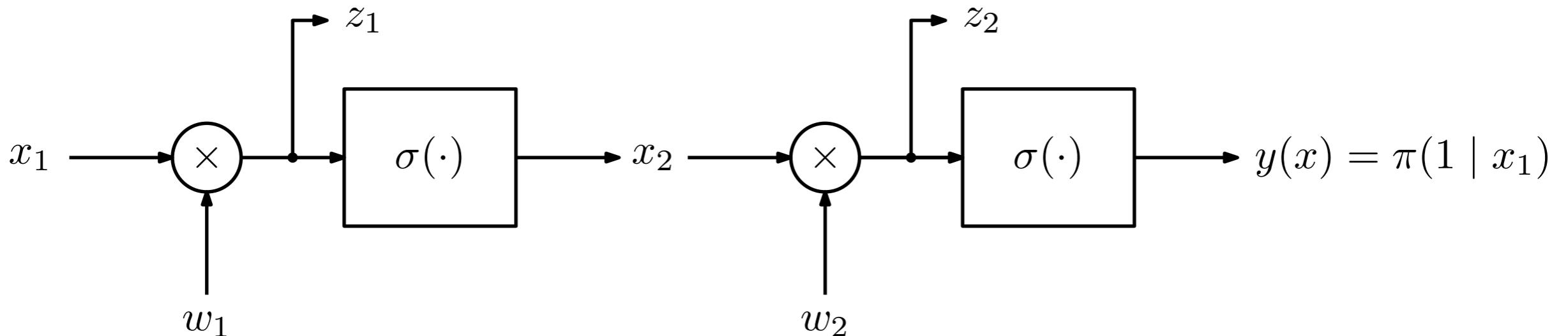
$$\frac{\partial \hat{L}_N}{\partial w_1} = \boxed{(y - a)} w_2 x_2 (1 - x_2) x_1$$

$\delta_2$

$$\frac{\partial \hat{L}_N}{\partial w_2} = (y - a) x_2$$

We can compute  
 $\delta_1$  from  $\delta_2$

# World's simplest net



$$\frac{\partial \hat{L}_N}{\partial w_1} = \boxed{\delta_1 x_1}$$

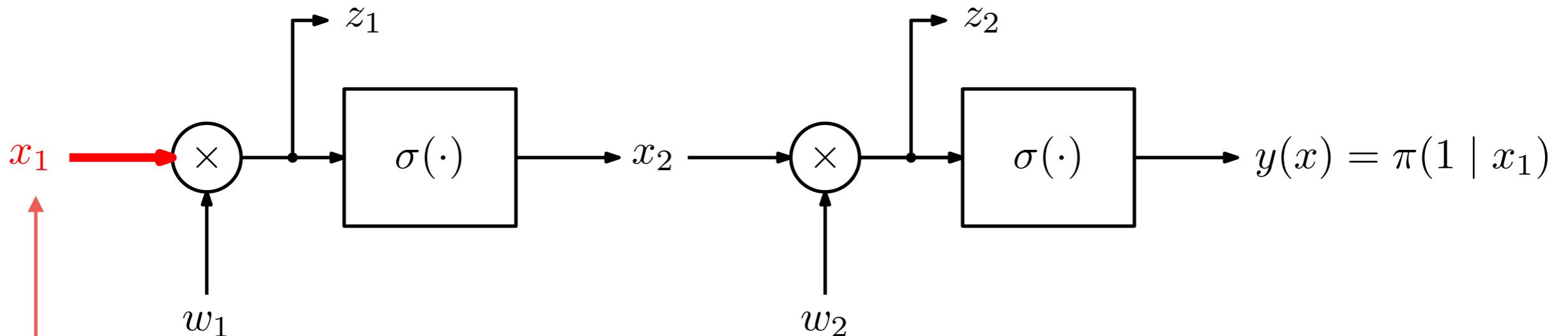
$$\frac{\partial \hat{L}_N}{\partial w_2} = \boxed{\delta_2 x_2}$$

$$\boxed{\delta_i = \frac{\partial \hat{L}_N}{\partial z_i}}$$

$$\boxed{\delta_{i-1} = \delta_i w_i \sigma'(x_i)}$$

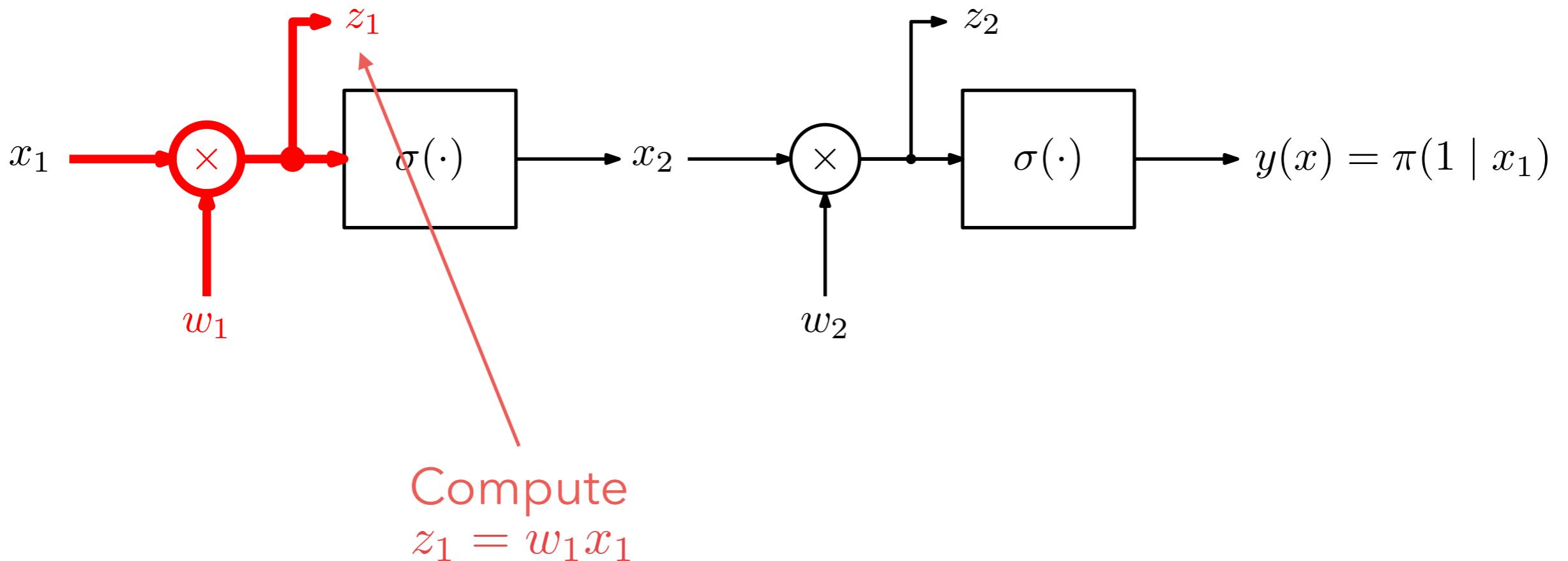
Let's repeat, now with pictures...

# Forward propagation

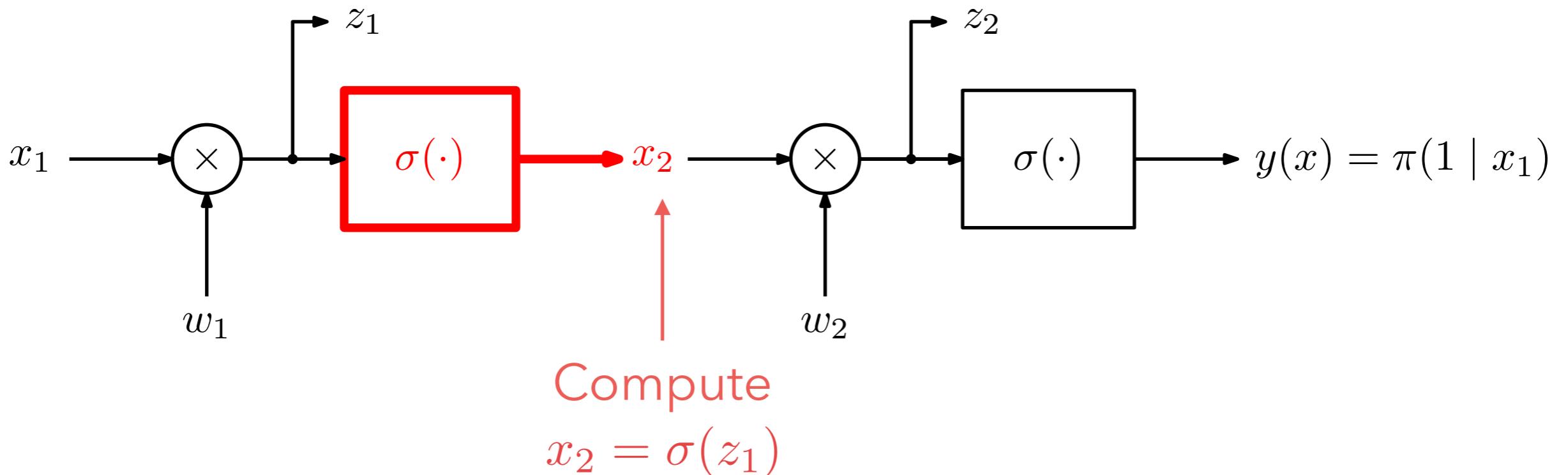


Given

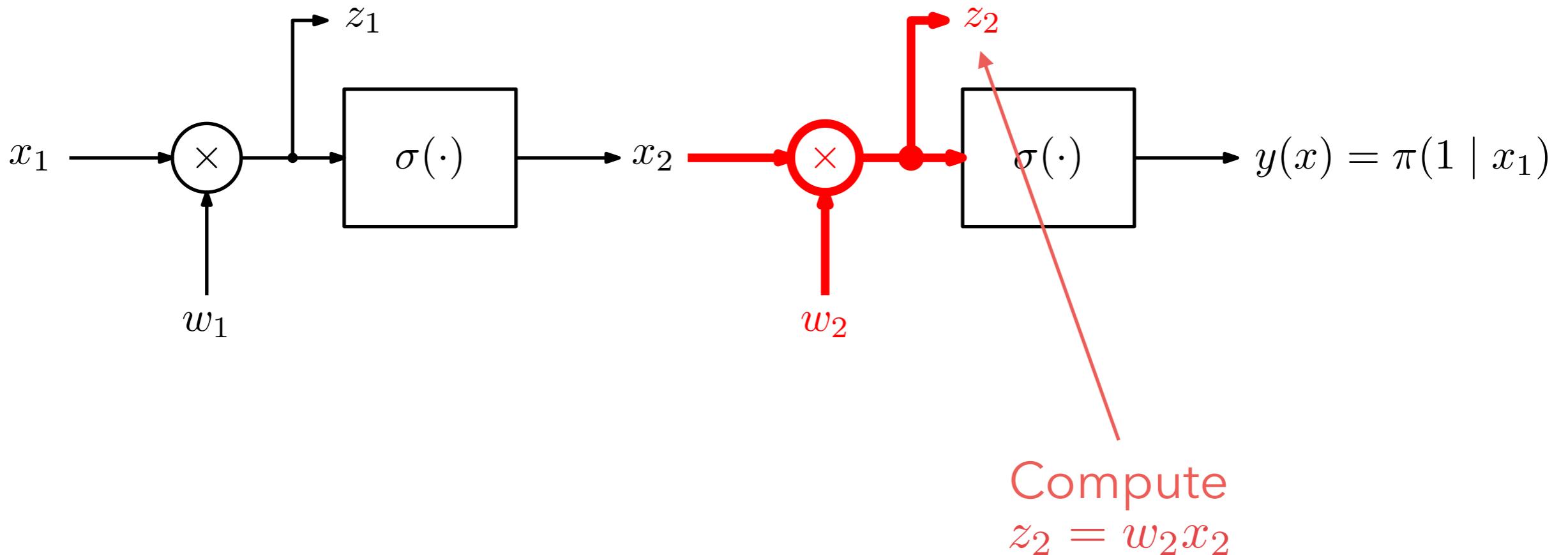
# Forward propagation



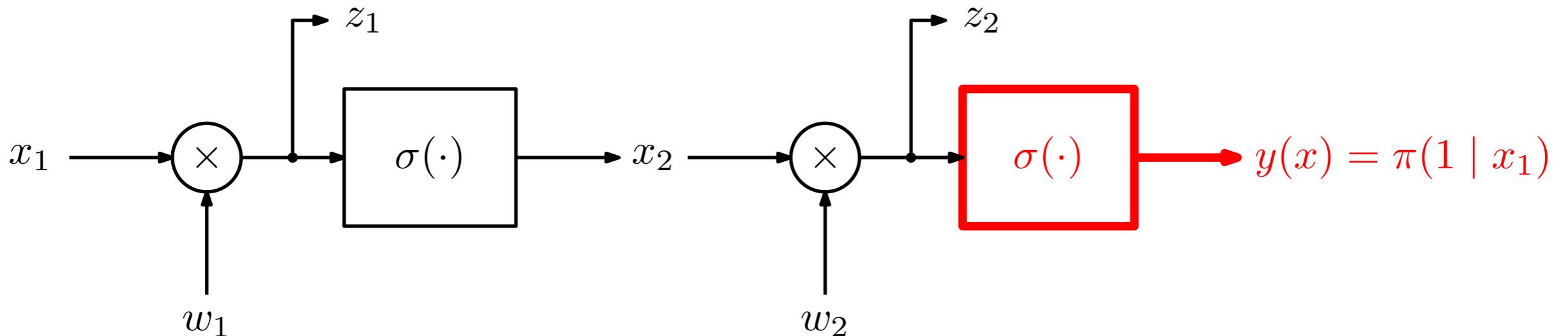
# Forward propagation



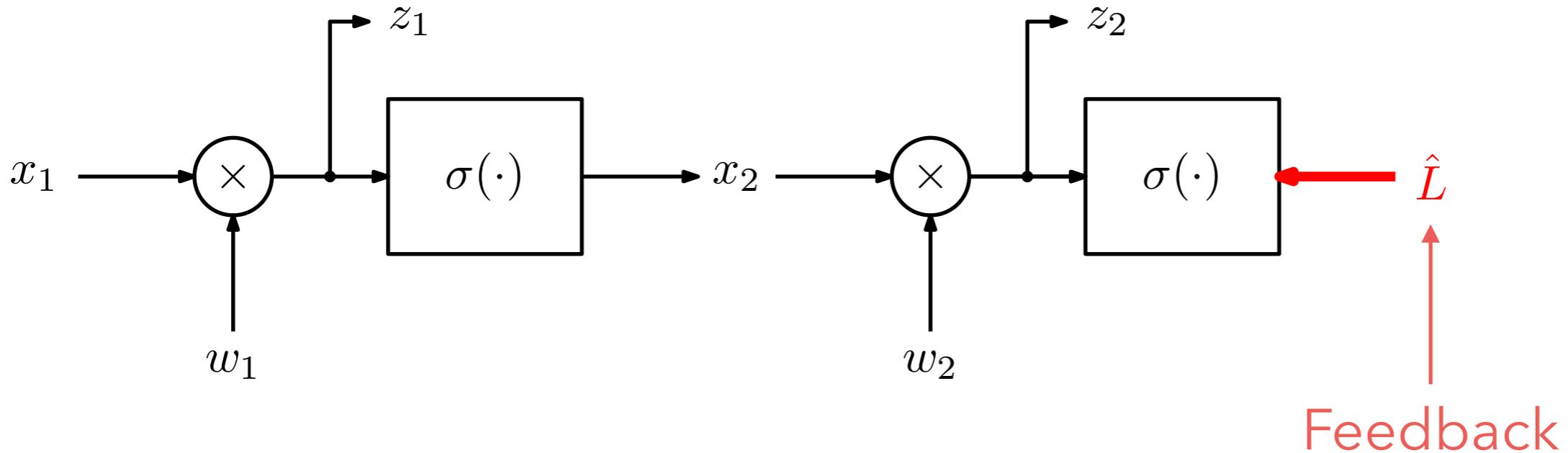
# Forward propagation



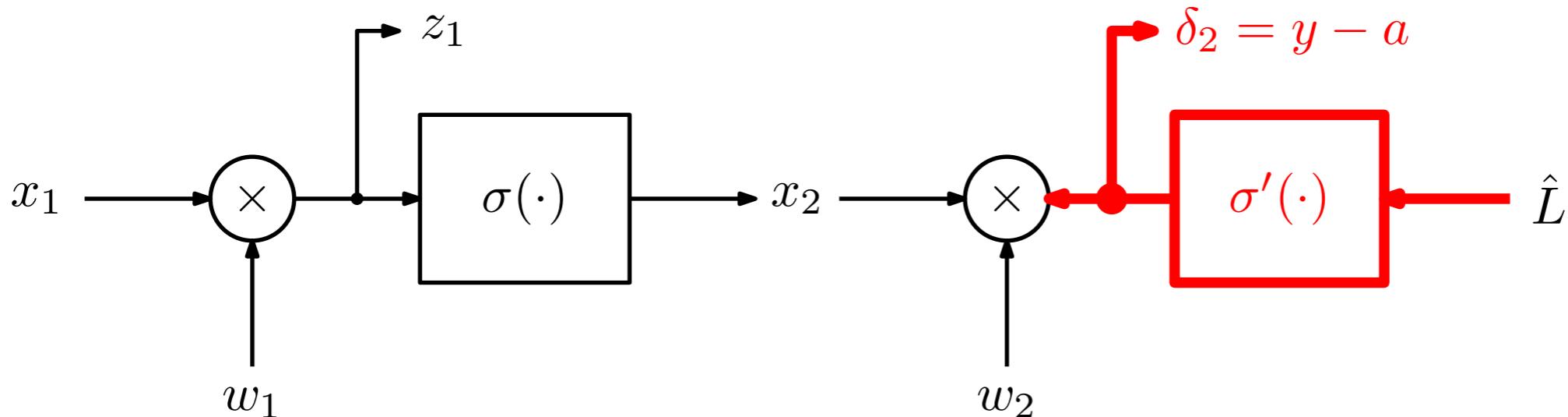
# Forward propagation



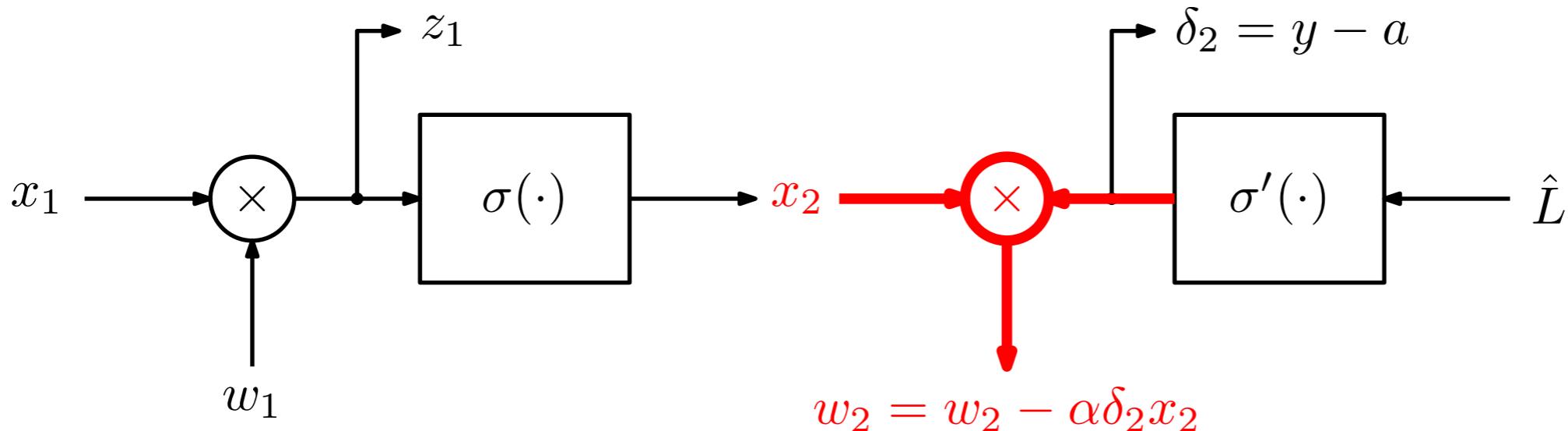
# Backward propagation



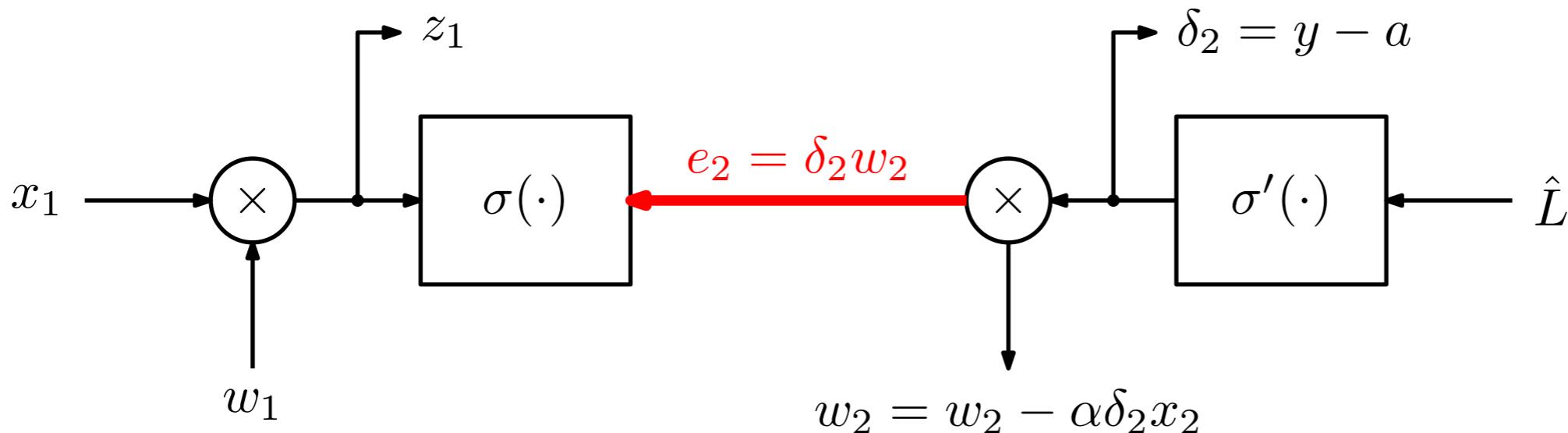
# Backward propagation



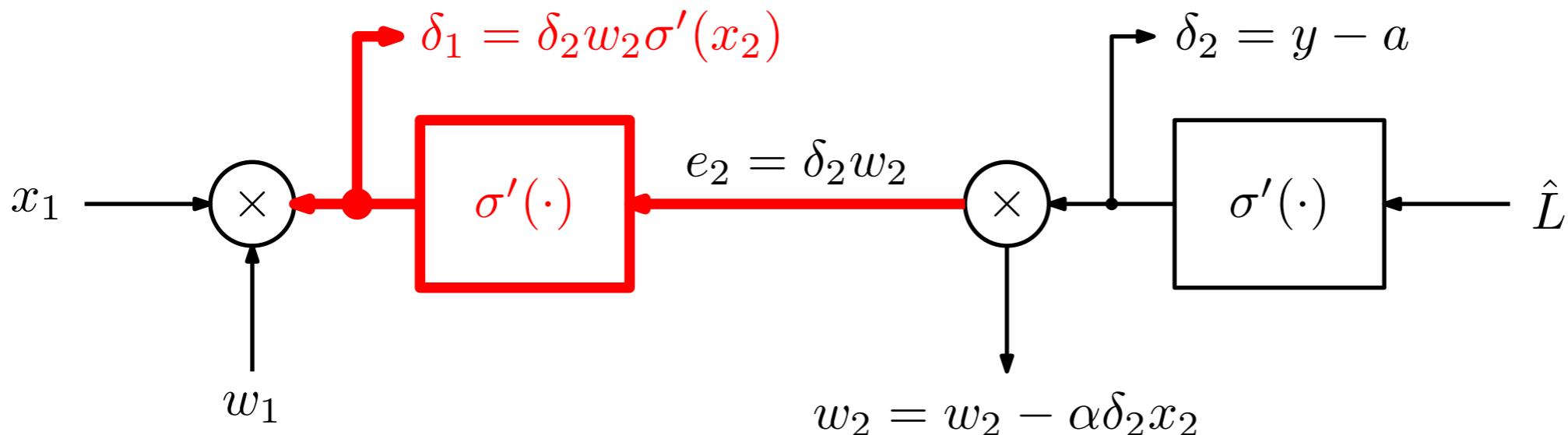
# Backward propagation



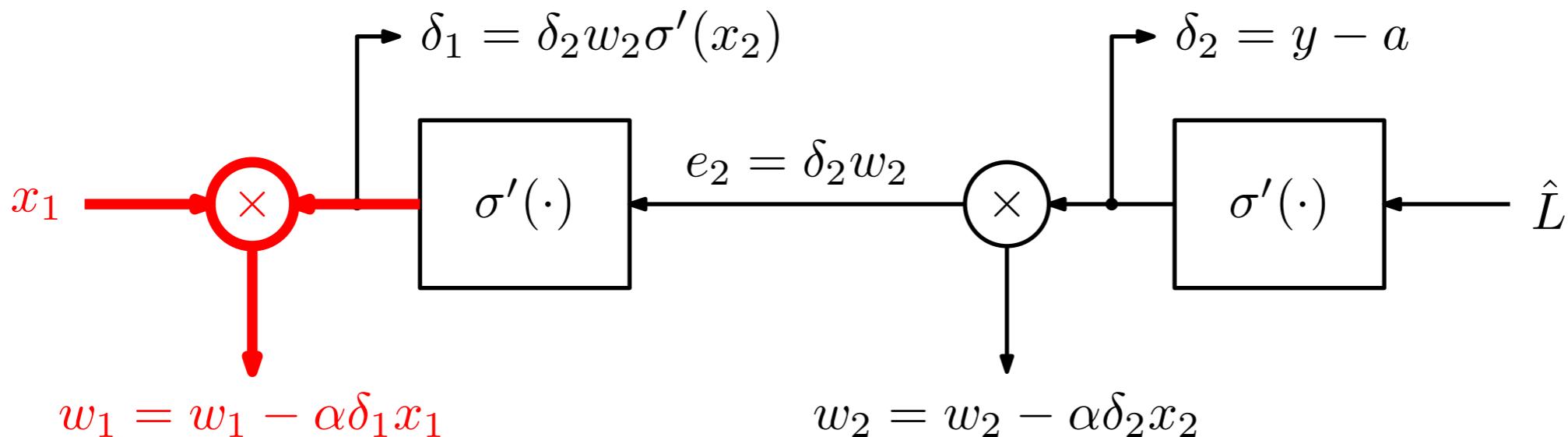
# Backward propagation



# Backward propagation



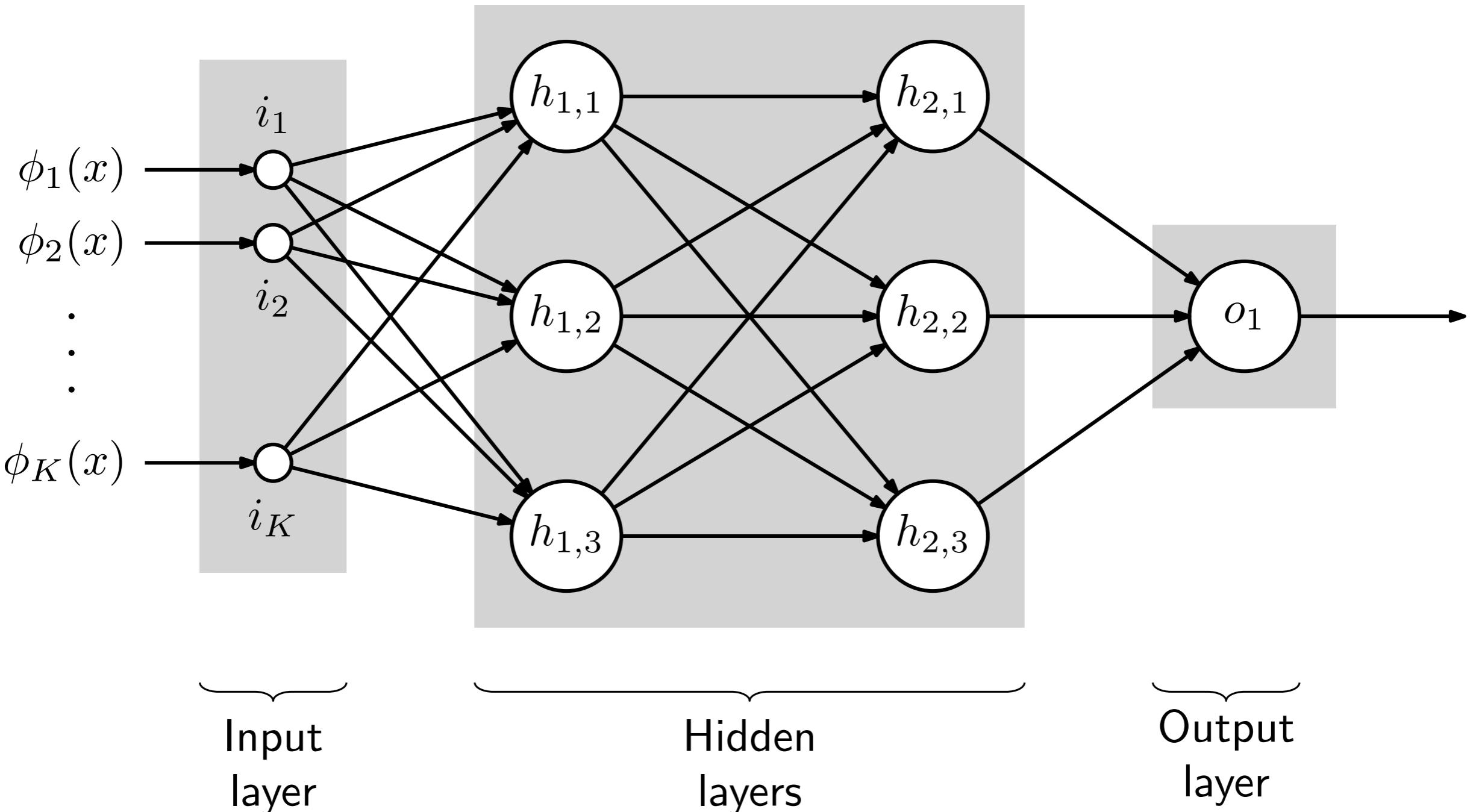
# Backward propagation



# Back-propagation

- Back-propagation is the standard approach to compute the gradient used in the training of neural networks
- Generalizes in a simple manner for more complicated networks

# Multilayer perceptron



# Multi-layer perceptron

- MLP can have multiple hidden layers
- Hidden layers and output layers can have varying number of neurons
- Layers are fully connected (all neurons in a layer connect to all neurons in the next layer)

# Back propagation for MLP

1. Initialize weights  $w_{ij}$  (connecting output of neuron i to neuron j) to small random values

2. Apply a data-point  $\mathbf{x}^*$  to the input layer

3. Propagate the signal through the network

$$y_i = \sigma(z_i) = \sigma(\mathbf{w}_{ji} \mathbf{y}_j)$$

4. Back-propagate  $\delta_i$  back through the network

$$\delta_i = \begin{cases} (a_i - a) & \text{if } i \text{ is an output unit} \\ \sigma'(z_i) \sum_k w_{ik} \delta_k & \text{otherwise} \end{cases}$$

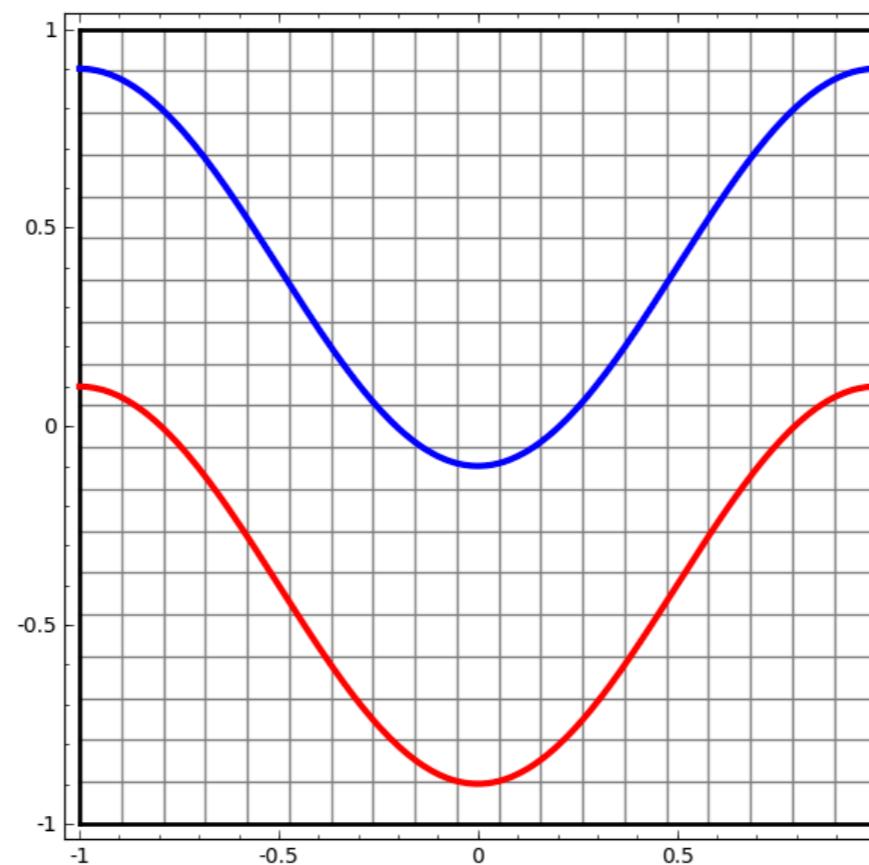
5. Update weights

# How do NN work?

- Output is like logistic regression
  - What do hidden layers do?

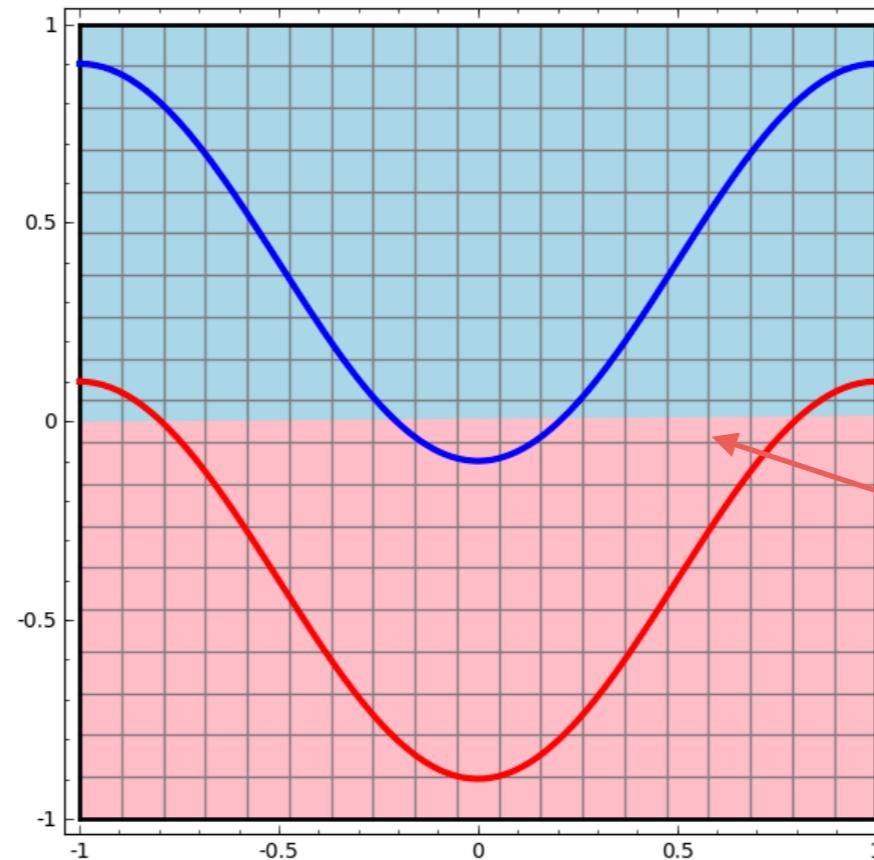
# How do NN work?

- Simple example:
  - Classify points as belonging to one of the two curves



# How do NN work?

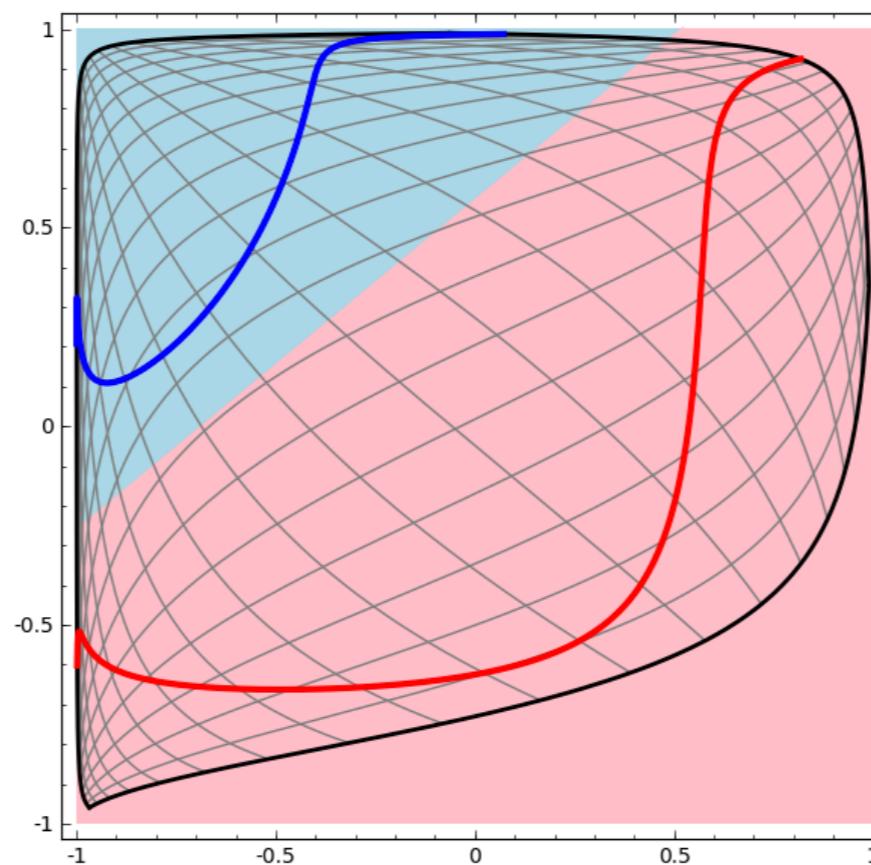
- Simple example:
  - Classify points as belonging to one of the two curves



A linear classifier (LR)  
does not work...

# How do NN work?

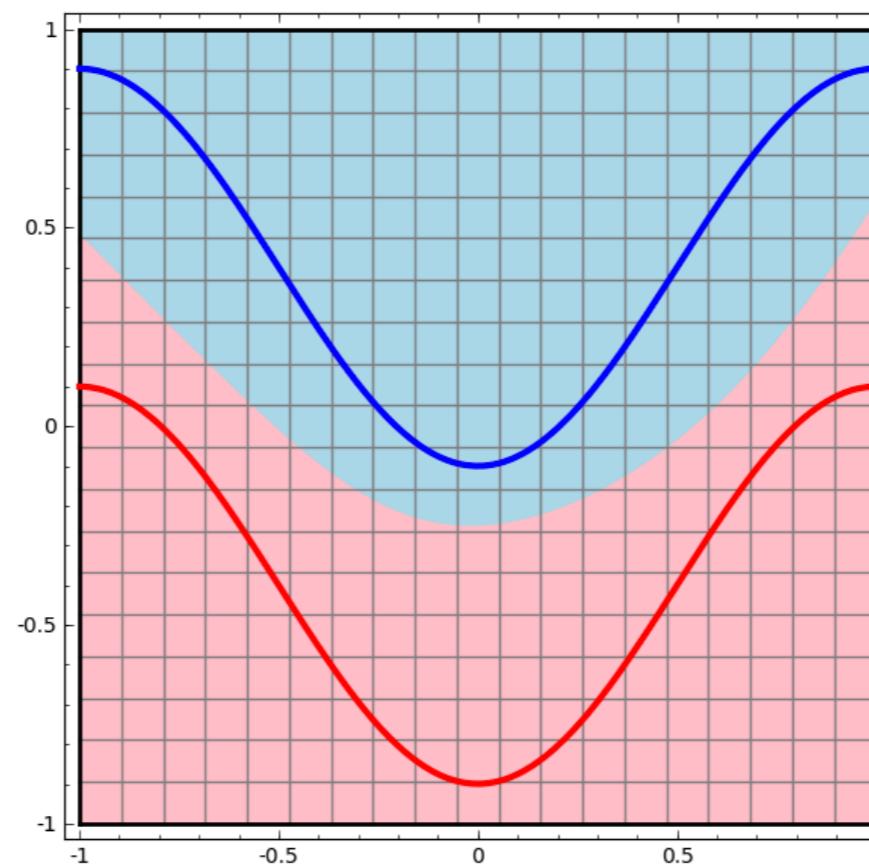
- Simple example:
  - Classify points as belonging to one of the two curves



Hidden layers  
“tweak” the space  
making the data  
linearly separable

# How do NN work?

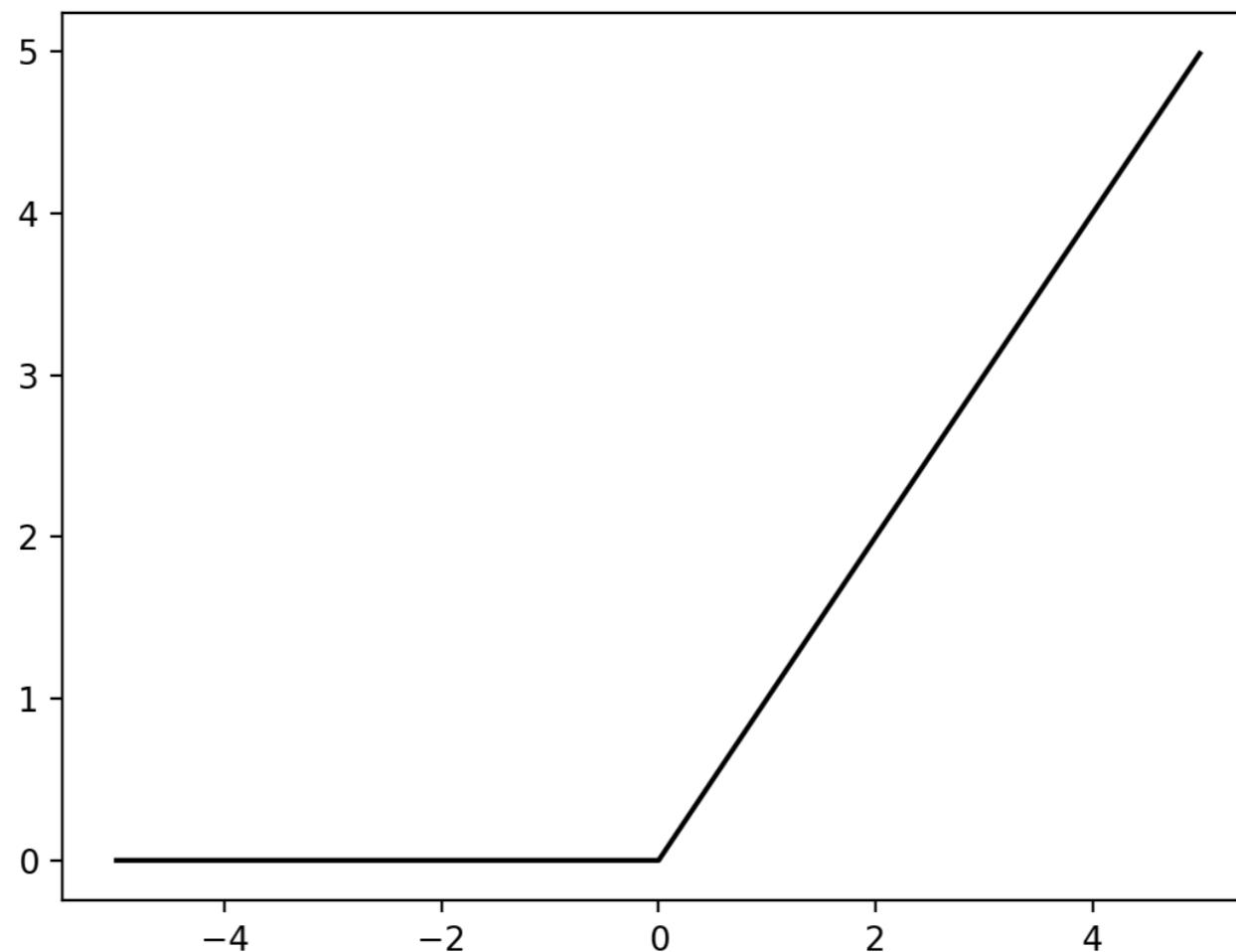
- Simple example:
  - Classify points as belonging to one of the two curves



In the original space,  
this is what the  
decision boundary  
looks like

# Final remarks

- Current practice has replaced the sigmoid function by a rectified linear unit (ReLU)



# Final remarks

- Current practice has replaced the sigmoid function by a rectified linear unit (ReLU)
- Neural networks are at the core of **deep learning**
- Neural networks used in deep learning are usually not MLP, although the structure near the output layers resembles it
- The training method is still back-propagation