



Acknowledgement: Many slides are adopted from Yaser Abu-Mostafa

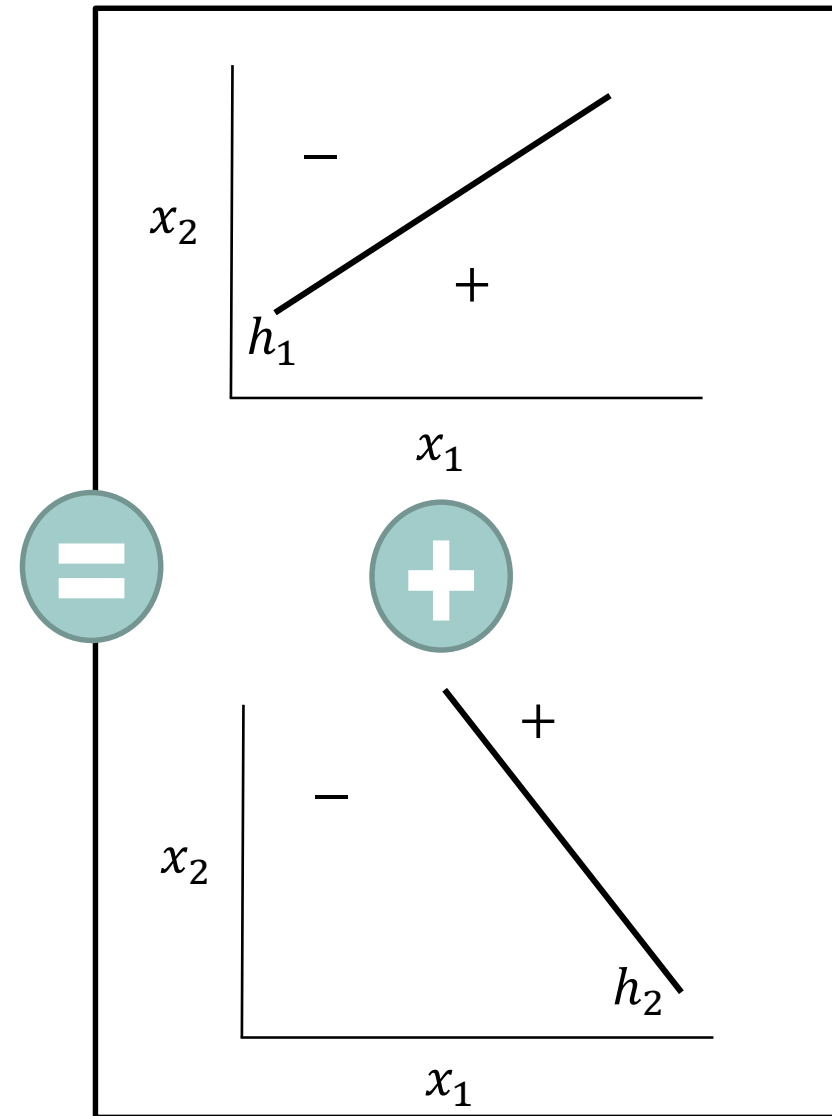
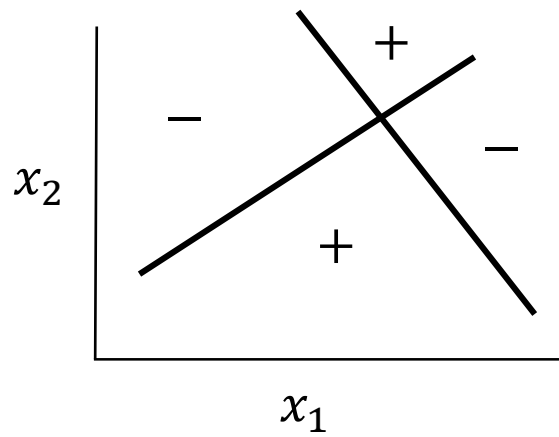


# Agenda

- Introduction
- Learning problem & linear classification
- Linear models: regression & logistic regression
- Non-linear transformation, overfitting & regularization
- Support Vector Machines and kernel learning
- **Neural Networks: shallow [and deep]**
- Theoretical foundation of supervised learning
- Unsupervised learning

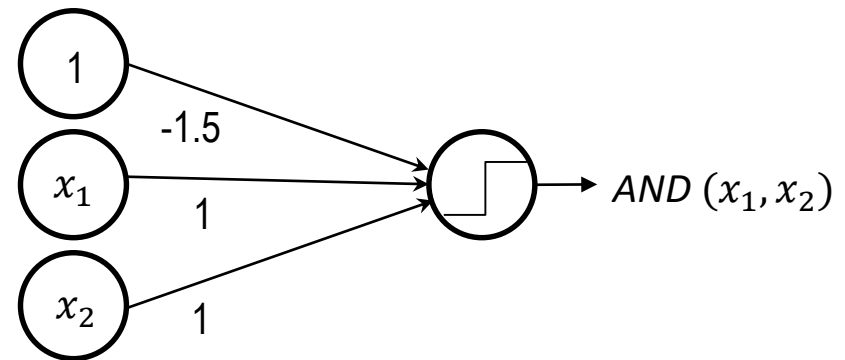
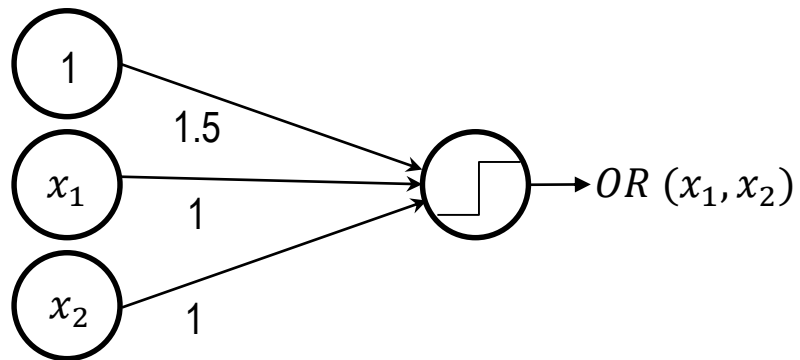


# The idea behind neural nets



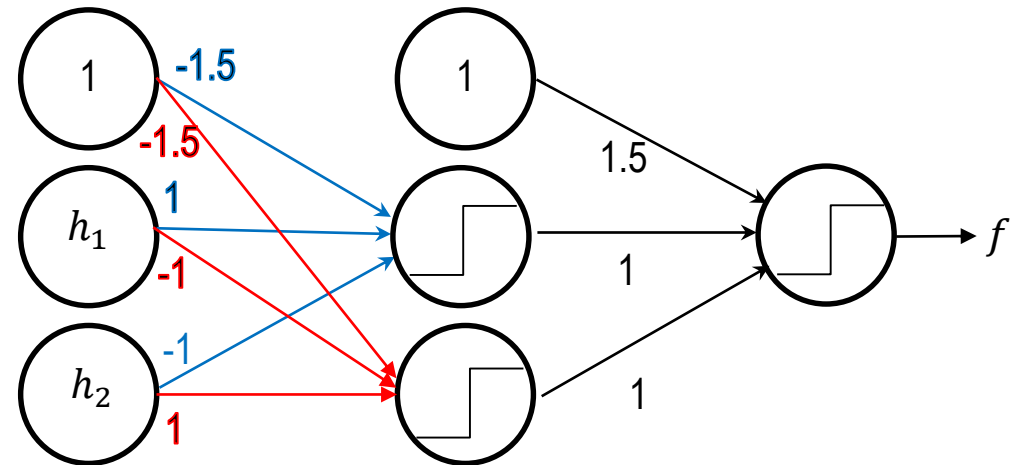
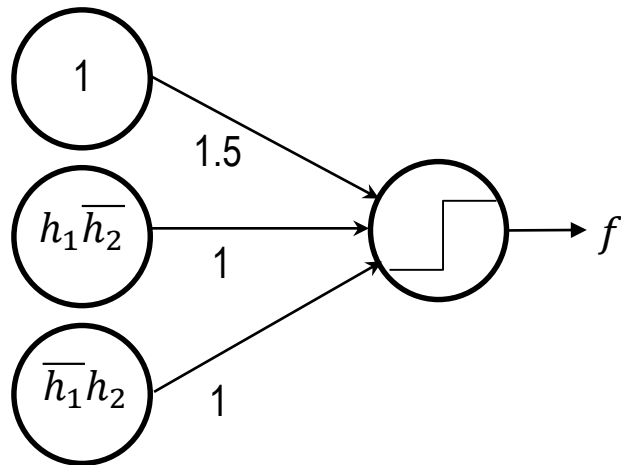


# Combining perceptrons to realize Boolean operators



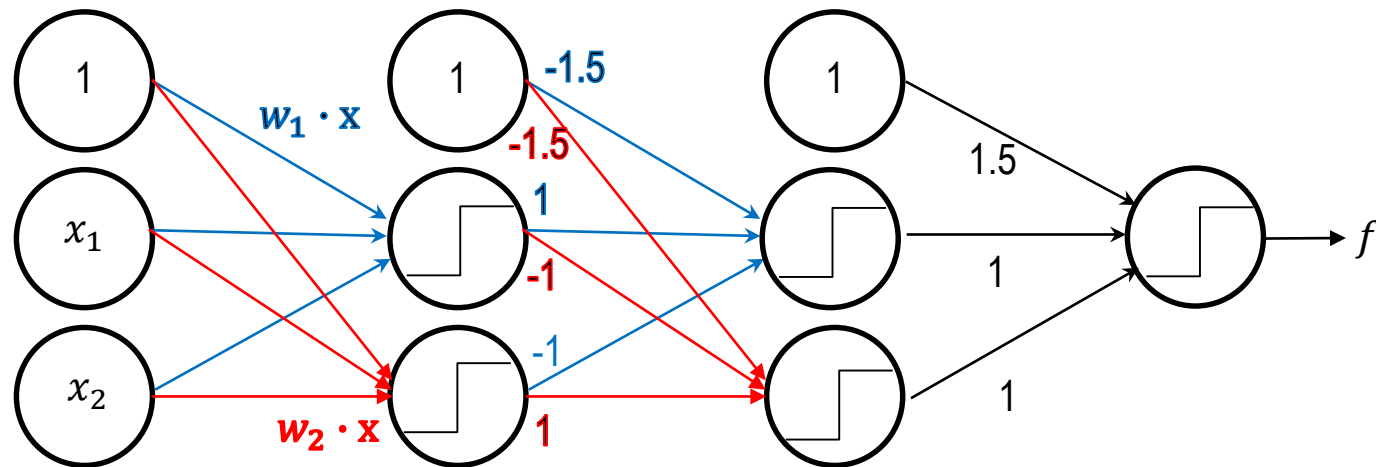


# Creating layers



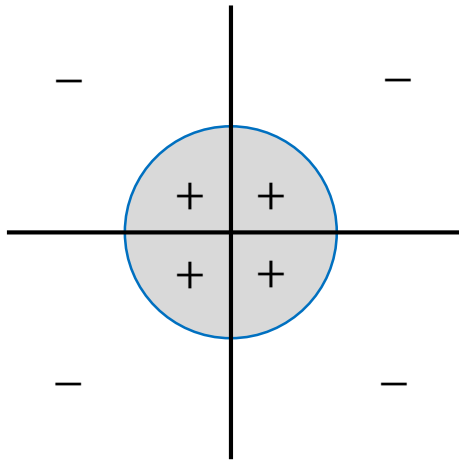


# The multilayer perceptron

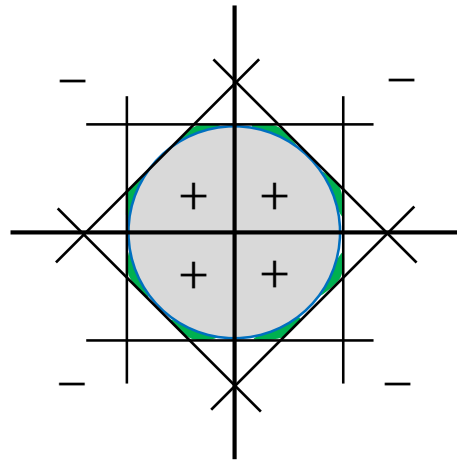




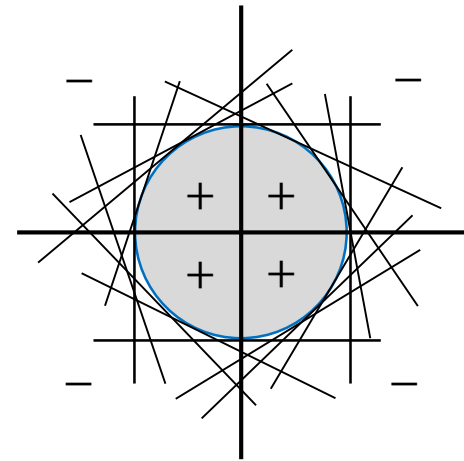
# A powerful model



Target



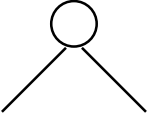
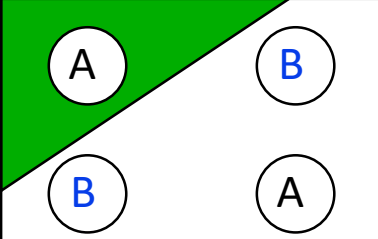
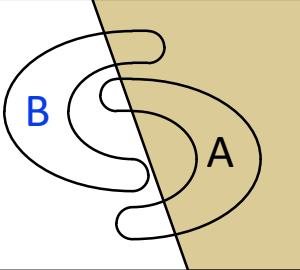
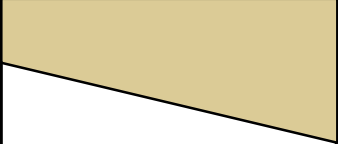
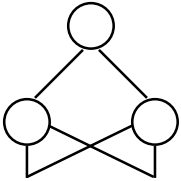
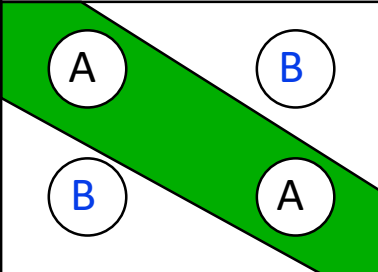
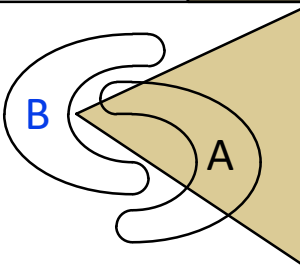
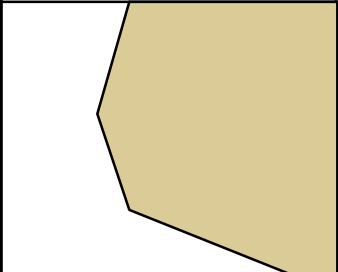
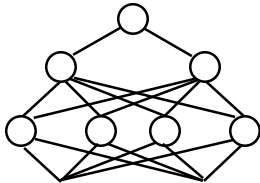
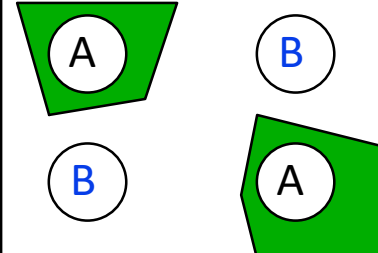
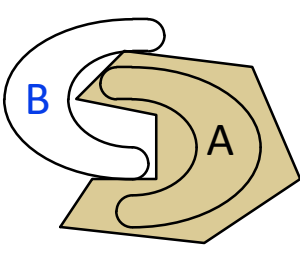
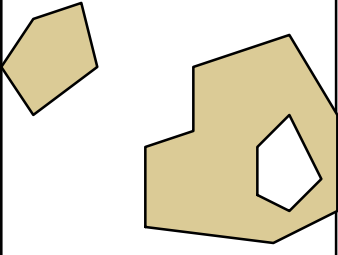
8 perceptrons



16 perceptrons



# The impact of layers

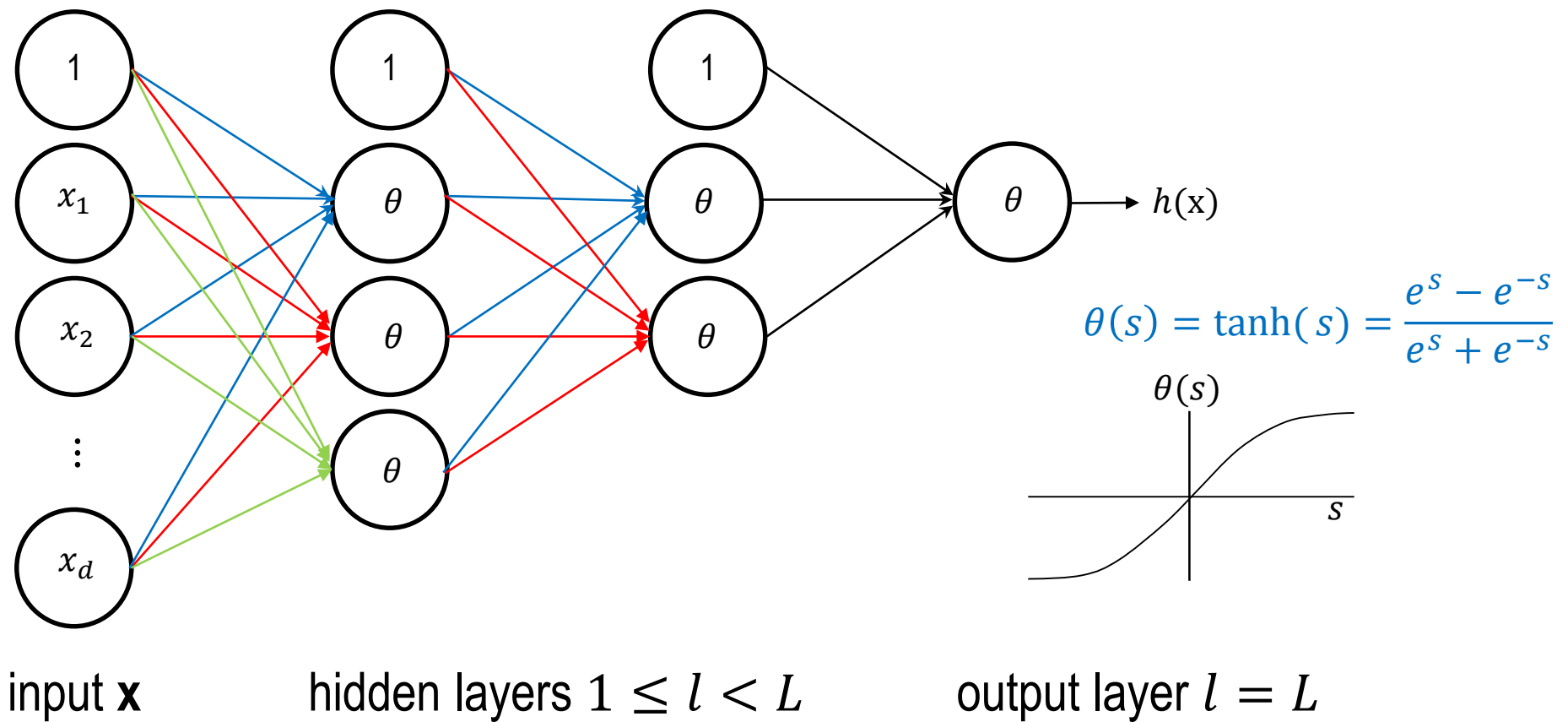
Structure	Types of Decision Regions	Exclusive-OR Problem	Classes with Meshed regions	Most General Region Shapes
Single-Layer 	Half Plane Bounded By Hyperplane			
Two-Layer 	Convex Open Or Closed Regions			
Three-Layer 	Arbitrary (Complexity Limited by No. of Nodes)			

Source: Mark Hoogendoorn (VU Amsterdam)



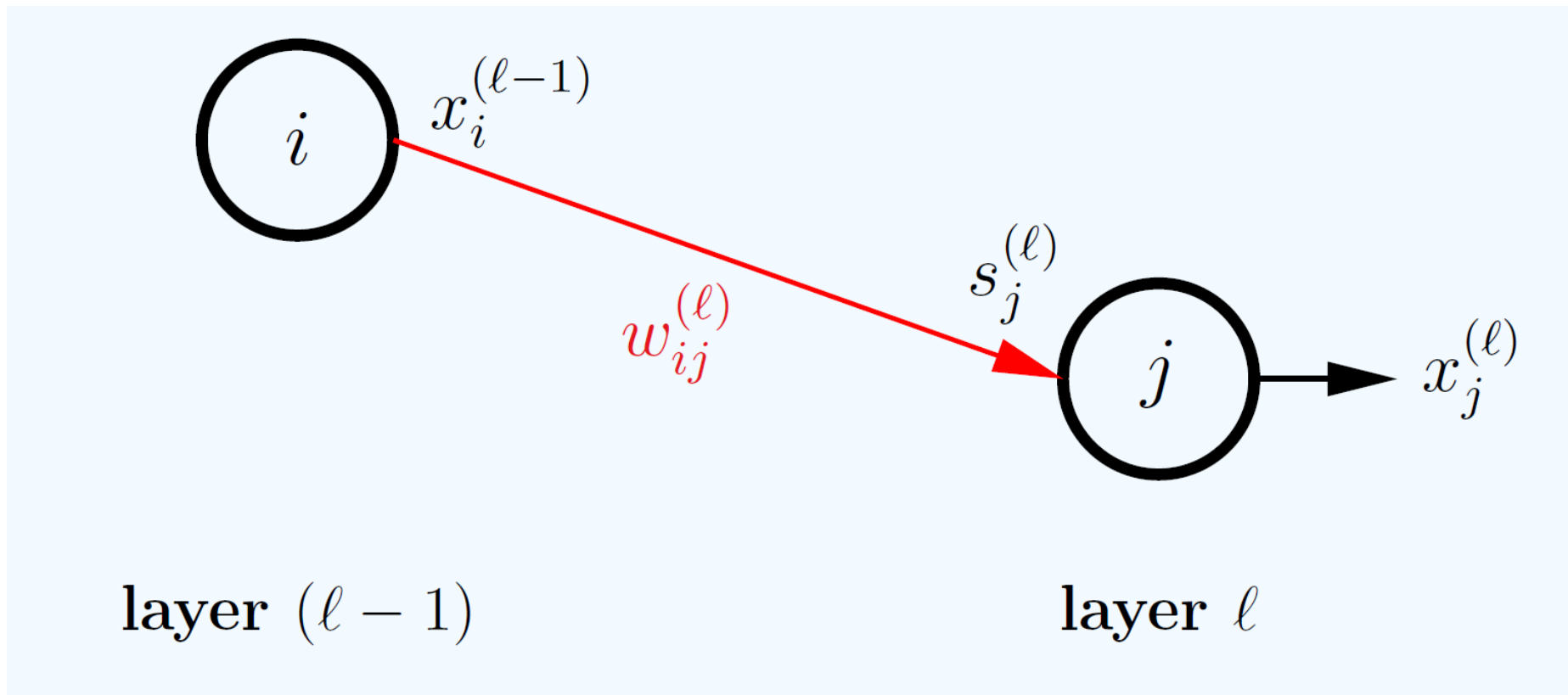


# A simple neural network





## The weight between two nodes



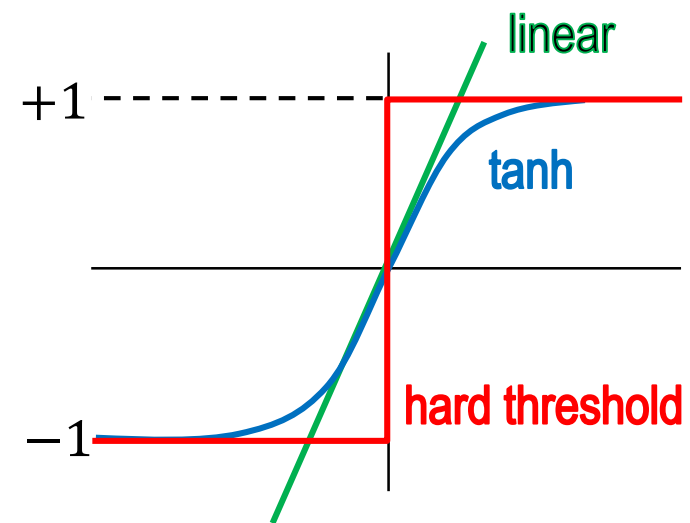


## Mathematical form

$$w_{ij}^{(l)} \begin{cases} 1 \leq l \leq L & \text{layers} \\ 0 \leq i \leq d^{(l-1)} & \text{inputs} \\ 1 \leq j \leq d^{(l)} & \text{outputs} \end{cases}$$

$$x_j^{(l)} = \theta \left( s_j^{(l)} \right) = \theta \left( \sum_{i=0}^{d^{(l-1)}} w_{ij}^{(l)} x_i^{(l-1)} \right)$$

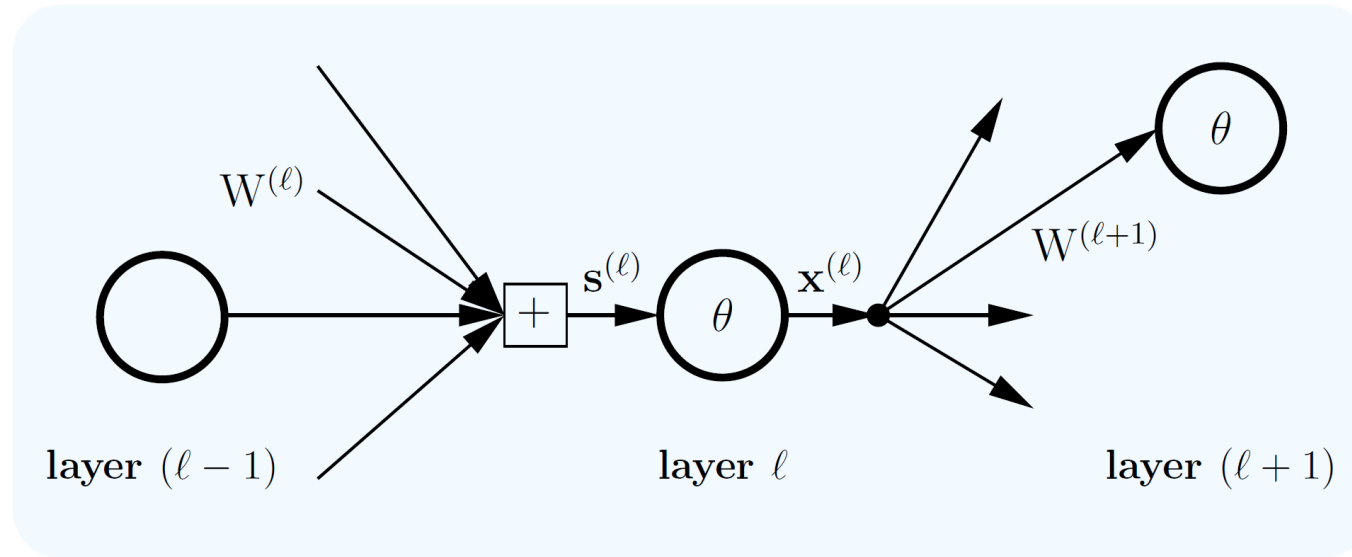
Apply  $\mathbf{x}$  to  $x_1^{(0)} \dots x_{d^{(0)}}^{(0)} \rightarrow \dots \rightarrow x_1^L = h(\mathbf{x})$



$$\theta(s) = \tanh(s) = \frac{e^s - e^{-s}}{e^s + e^{-s}}$$



## Vectorized version



layer  $\ell$  parameters

signals in	$\mathbf{s}^{(\ell)}$	$d^{(\ell)}$ dimensional input vector
outputs	$\mathbf{x}^{(\ell)}$	$d^{(\ell)} + 1$ dimensional output vector
weights in	$W^{(\ell)}$	$(d^{(\ell-1)} + 1) \times d^{(\ell)}$ dimensional matrix
weights out	$W^{(\ell+1)}$	$(d^{(\ell)} + 1) \times d^{(\ell+1)}$ dimensional matrix



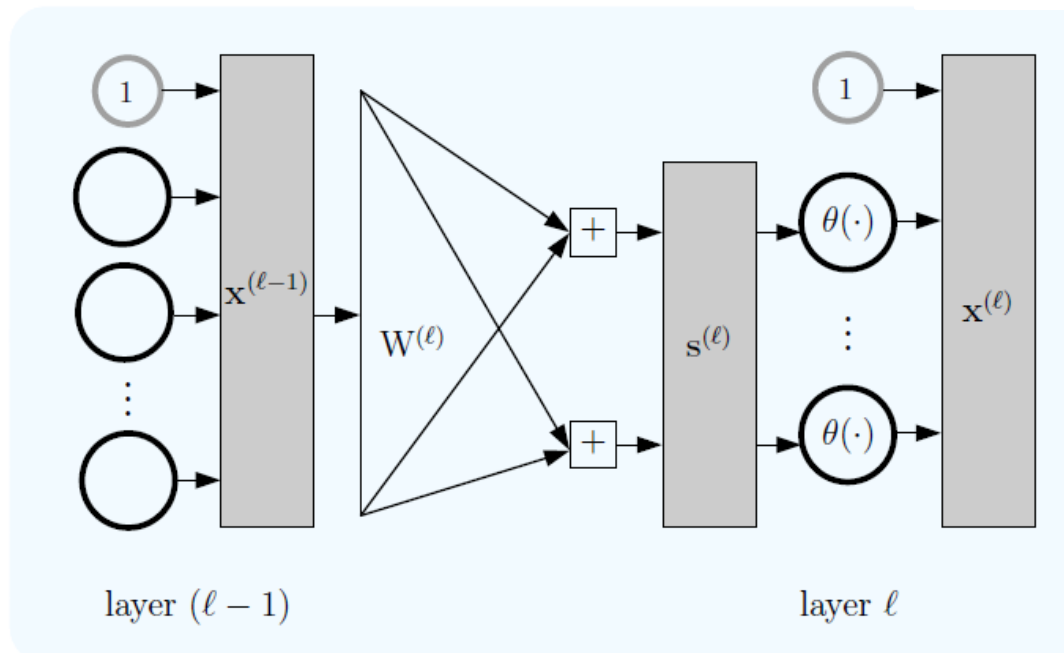
# Forward Propagation

- 1:  $\mathbf{x}^{(0)} \leftarrow \mathbf{x}$
- 2: **for**  $\ell = 1$  **to**  $L$  **do**
- 3:    $\mathbf{s}^{(\ell)} \leftarrow (\mathbf{W}^{(\ell)})^\top \mathbf{x}^{(\ell-1)}$
- 4:    $\mathbf{x}^{(\ell)} \leftarrow \begin{bmatrix} 1 \\ \theta(\mathbf{s}^{(\ell)}) \end{bmatrix}$
- 5:  $h(\mathbf{x}) = \mathbf{x}^{(L)}$

[Initialization]

[Forward Propagation]

[Output]





# Applying stochastic gradient descent

All the weights  $\mathbf{w} = \{w_{ij}^{(l)}\}$  determine  $h(\mathbf{x})$

Error on example  $(x_n, y_n)$  is

$$e(h(x_n), y_n) = e(\mathbf{w})$$

To implement SGD, we need the gradient

$$\nabla e(\mathbf{w}): \frac{\partial e(\mathbf{w})}{\partial w_{ij}^{(l)}} \quad \text{for all } i, j, l$$



# Remember the chain rule?



## Computing $\frac{\partial e(\mathbf{w})}{\partial w_{ij}^{(l)}}$

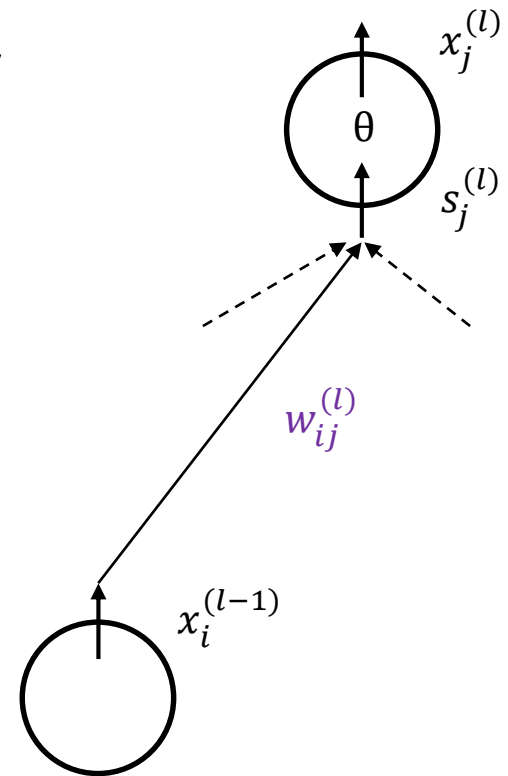
We can evaluate  $\frac{\partial e(\mathbf{w})}{\partial w_{ij}^{(l)}}$  one by one: analytically or numerically

A trick for efficient computation:

$$\frac{\partial e(\mathbf{w})}{\partial w_{ij}^{(l)}} = \frac{\partial e(\mathbf{w})}{\partial s_j^{(l)}} \times \frac{\partial s_j^{(l)}}{\partial w_{ij}^{(l)}}$$

We have  $\frac{\partial s_j^{(l)}}{\partial w_{ij}^{(l)}} = x_i^{(l-1)}$

We only need:  $\frac{\partial e(\mathbf{w})}{\partial s_j^{(l)}} = \delta_j^{(l)}$







## $\delta$ for the final layer

$$\delta_j^{(l)} = \frac{\partial e(\mathbf{w})}{\partial s_j^{(l)}}$$

For the final layer  $l = L$  and  $j = 1$  :

$$\delta_1^{(L)} = \frac{\partial e(\mathbf{w})}{\partial s_1^{(L)}}$$

$$e(\mathbf{w}) := (x_1^{(L)} - y_n)^2$$

$$x_1^{(L)} = \theta(s_1^{(L)})$$

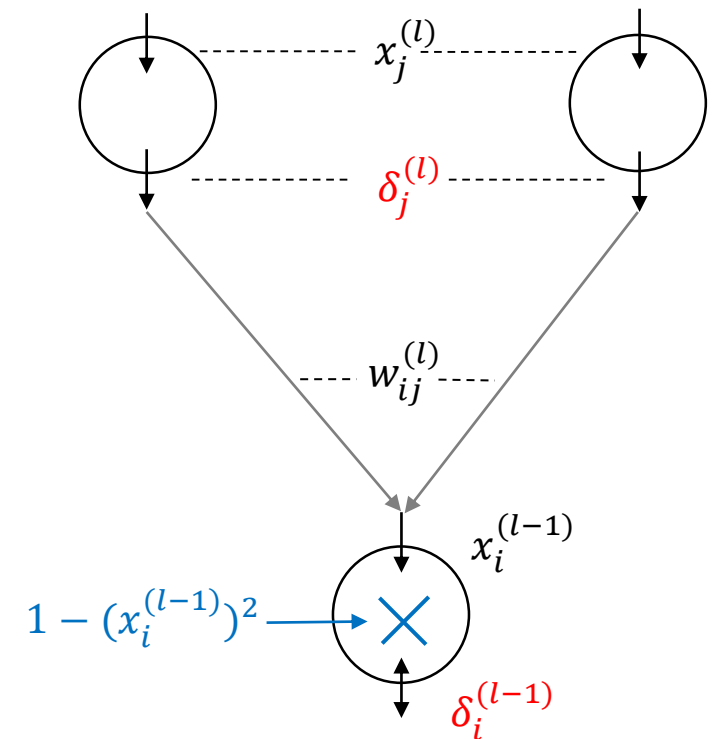
$$\theta'(s) = 1 - \theta^2(s) \quad \text{for the tanh}$$



## Back propagation $\delta$

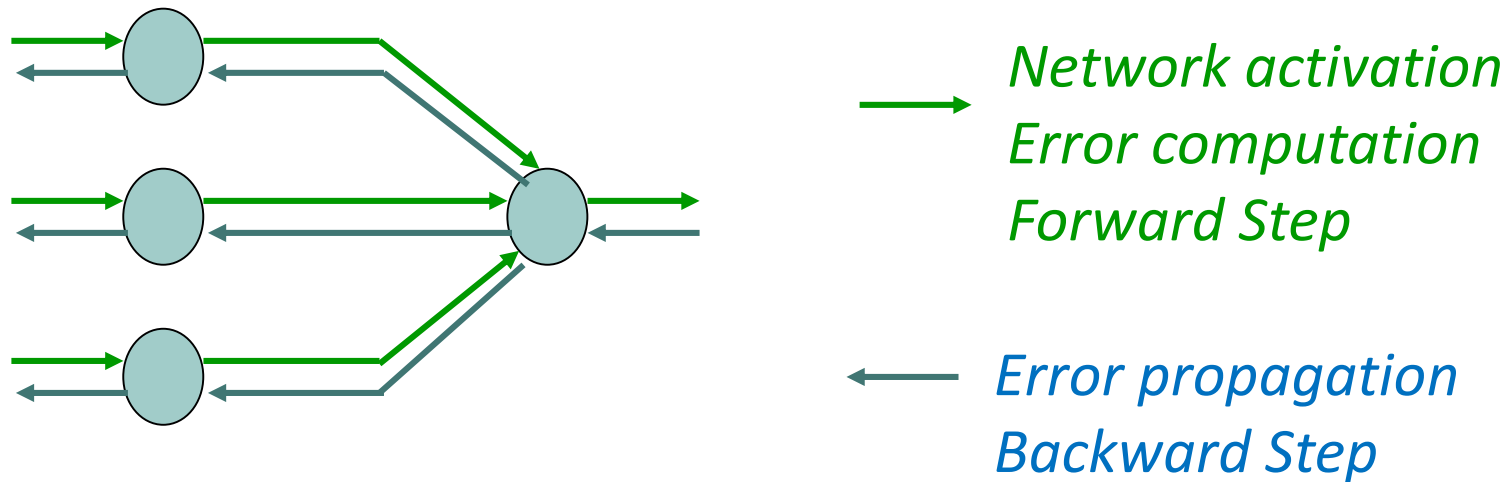
$$\begin{aligned}\delta_i^{(l-1)} &= \frac{\partial e(\mathbf{w})}{\partial s_i^{(l-1)}} \\ &= \sum_{j=1}^{d^{(l)}} \frac{\partial e(\mathbf{w})}{\partial s_j^{(l)}} \times \frac{\partial s_j^{(l)}}{\partial x_i^{(l-1)}} \times \frac{\partial x_i^{(l-1)}}{\partial s_i^{(l-1)}} \\ &= \sum_{j=1}^{d^{(l)}} \delta_j^{(l)} \times w_{ij}^{(l)} \times \theta'(s_i^{(l-1)})\end{aligned}$$

$$\delta_i^{(l-1)} = \left(1 - (x_i^{(l-1)})^2\right) \sum_{j=1}^{d^{(l)}} w_{ij}^{(l)} \delta_j^{(l)}$$





# Combing feedforward and backpropagation



Source: Mark Hoogendoorn (VU Amsterdam)



# Backpropagation algorithm

Initialize all weights  $w_{ij}^{(l)}$  **at random**

**for**  $t = 0, 1, 2, \dots$  **do**

Pick  $n \in \{1, 2, \dots, N\}$

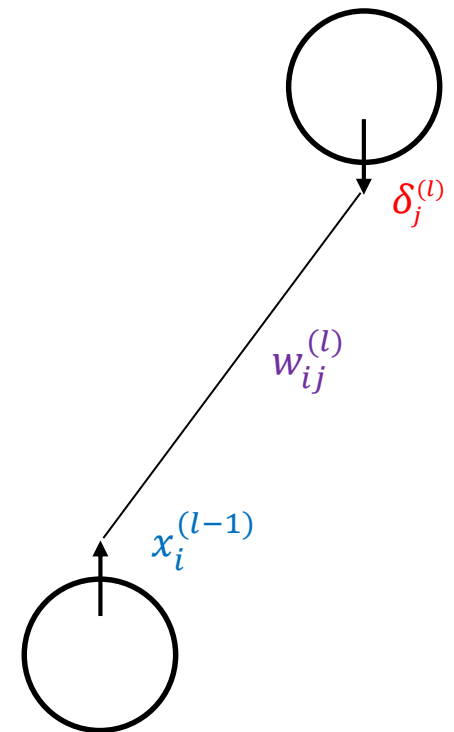
*Forward:* Compute all  $x_j^{(l)}$

*Backward:* Compute all  $\delta_j^{(l)}$

Update the weights:  $w_{ij}^{(l)} \leftarrow w_{ij}^{(l)} - \eta x_i^{(l-1)} \delta_j^{(l)}$

Iterate to the next step until it is time to stop

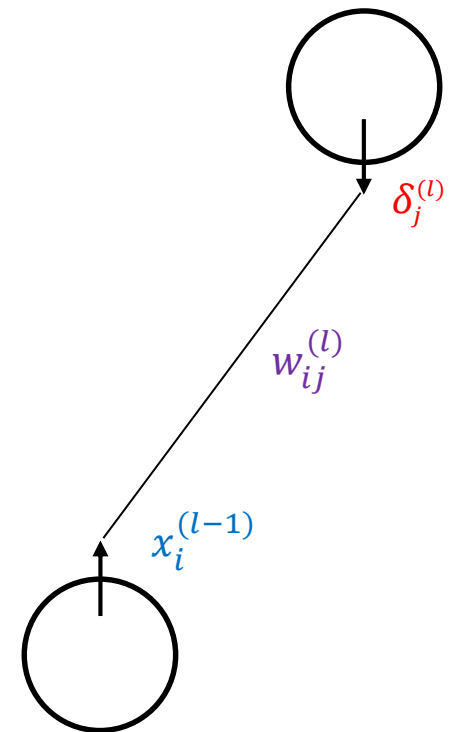
Return the final weights  $w_{ij}^{(l)}$





# A few words on the Backpropagation Algorithm

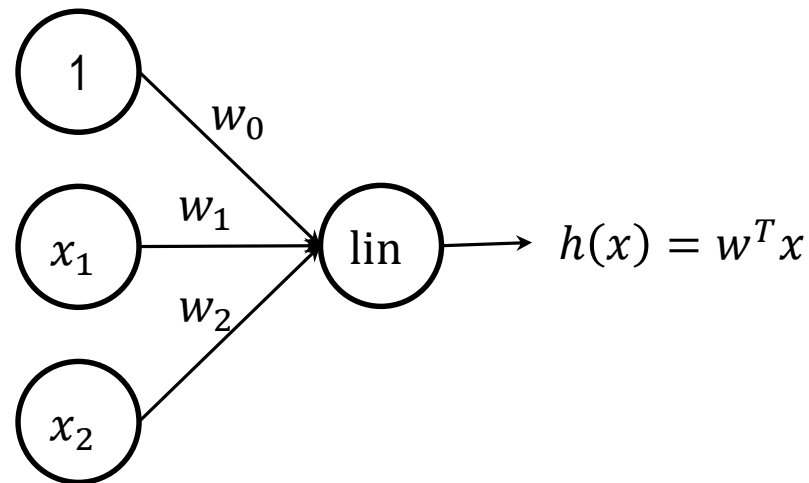
- Backpropagation Algorithm (BA) searches for weight values that minimize the total error of the network by iterative application of forward and backward pass
- Convergence is not guaranteed (e.g. learning rate too big or too small)
- If BA converges it is not guaranteed, that we reached the global minimum





## A final word: neural networks for regression

- The backprop algorithm does not heavily depend on the activation function
- If we substitute the  $\tanh(s)$  output node just by the linear function  $(s)$  something magic happens



- Any continuous function can be approximated by an NN with 2 hidden layers



# Deep Learning

