

Neural Networks 1: Multilayer Perceptron

Multiple layers
Universal Approximation
The Neural Network

The Neural Network - Biologically Inspired

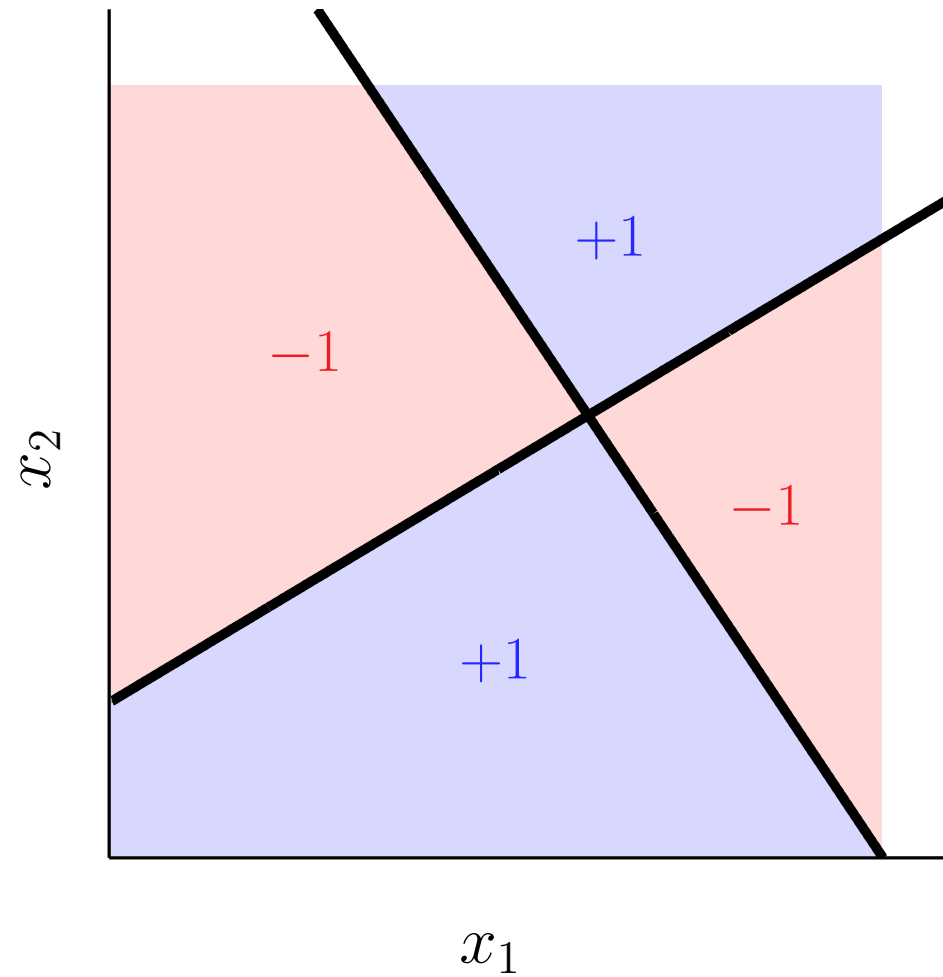


Planes Don't Flap Wings to Fly

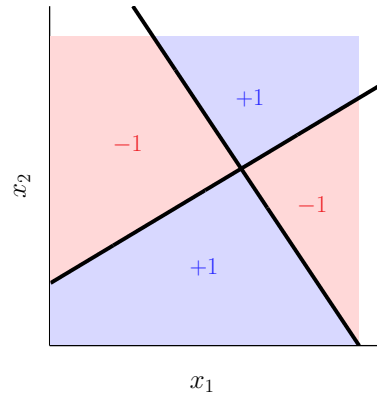


Engineering success may start with biological inspiration, but then take a totally different path.

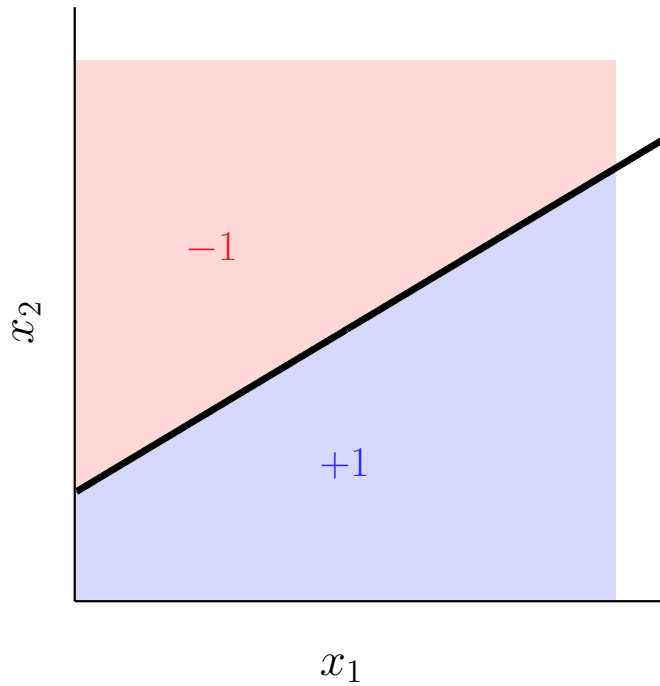
XOR: A Limitation of the Linear Model



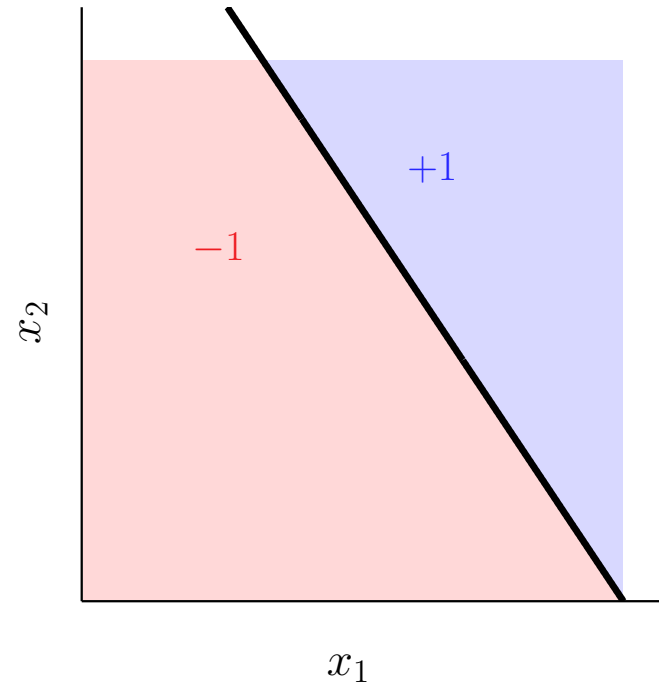
Decomposing XOR



$$f = h_1 \overline{h_2} + \overline{h_1} h_2$$



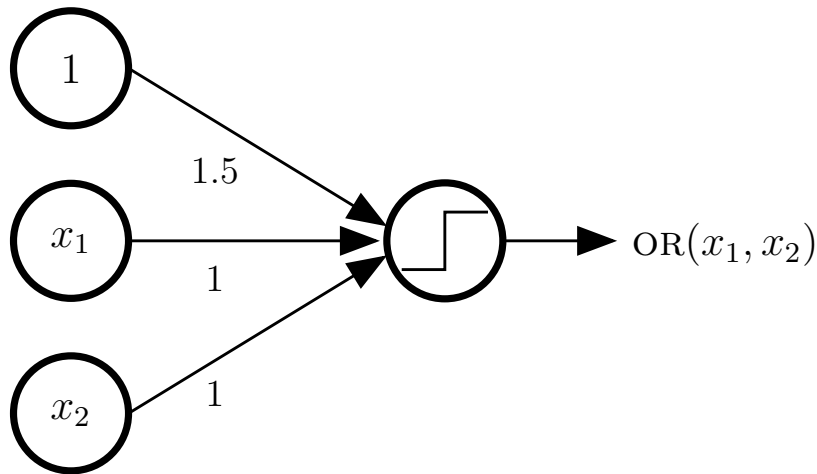
$$h_1(\mathbf{x}) = \text{sign}(\mathbf{w}_1^T \mathbf{x})$$



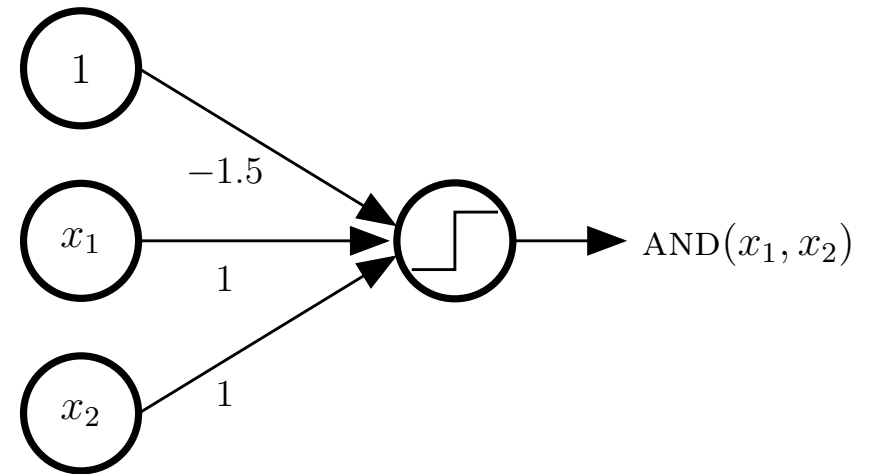
$$h_2(\mathbf{x}) = \text{sign}(\mathbf{w}_2^T \mathbf{x})$$

Perceptrons for OR and AND

$$\text{OR}(x_1, x_2) = \text{sign}(x_1 + x_2 + 1.5)$$

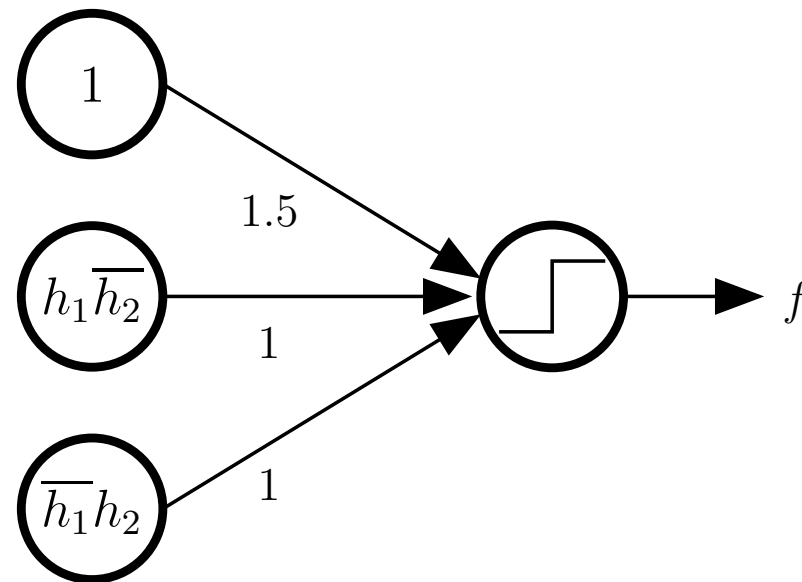


$$\text{AND}(x_1, x_2) = \text{sign}(x_1 + x_2 - 1.5)$$



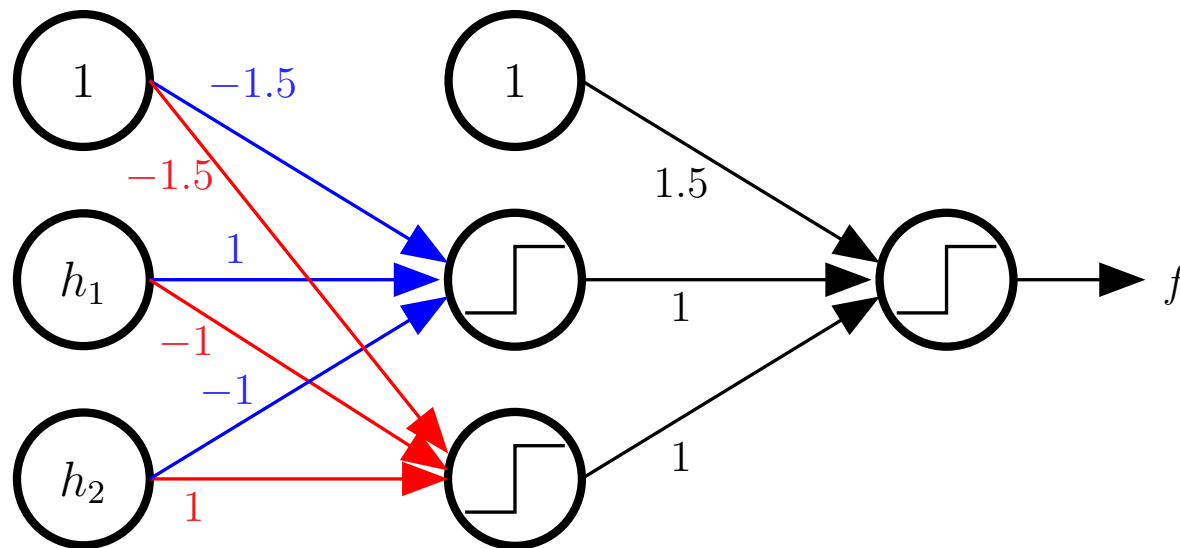
Representing f Using OR and AND

$$f = h_1 \overline{h_2} + \overline{h_1} h_2$$



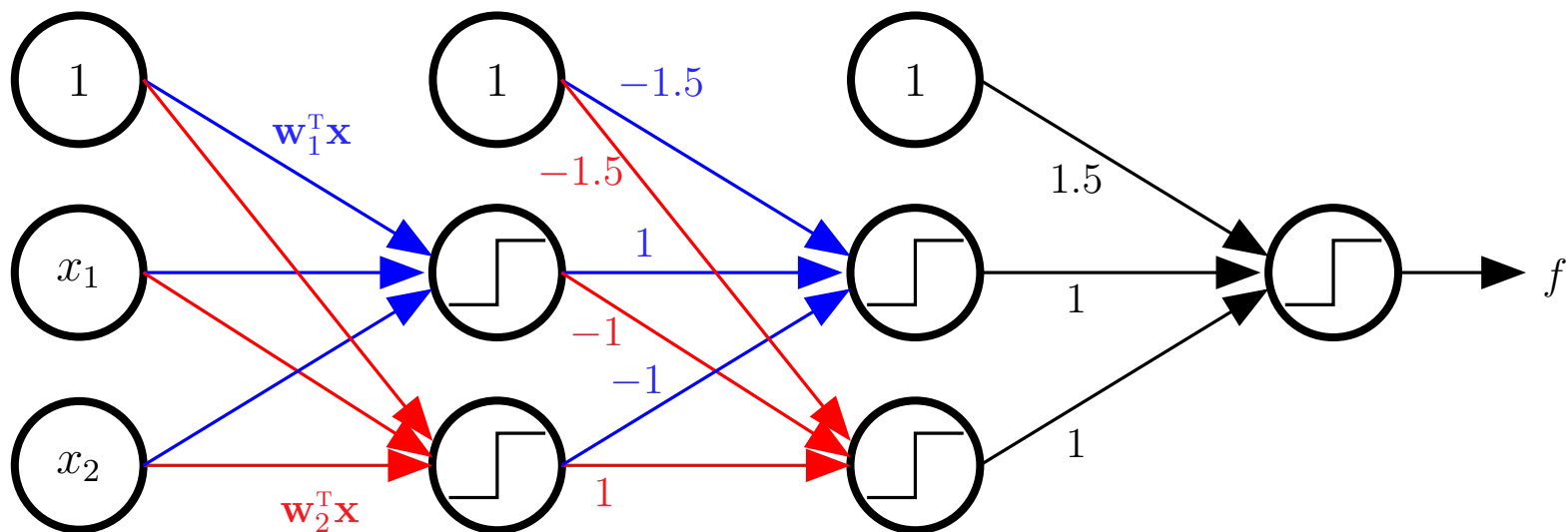
Representing f Using OR and AND

$$f = h_1 \overline{h_2} + \overline{h_1} h_2$$

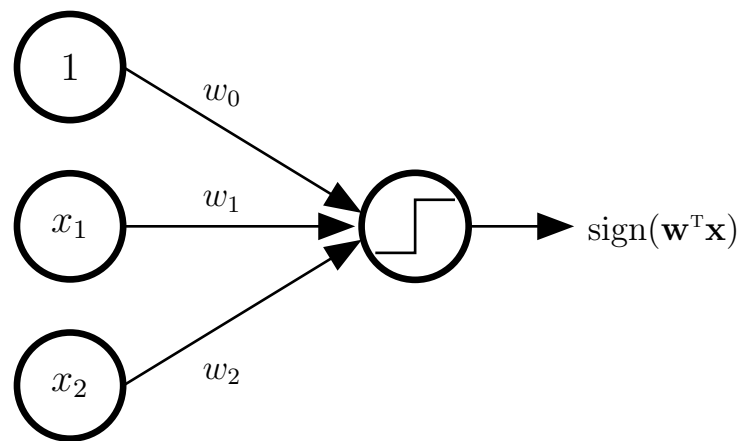
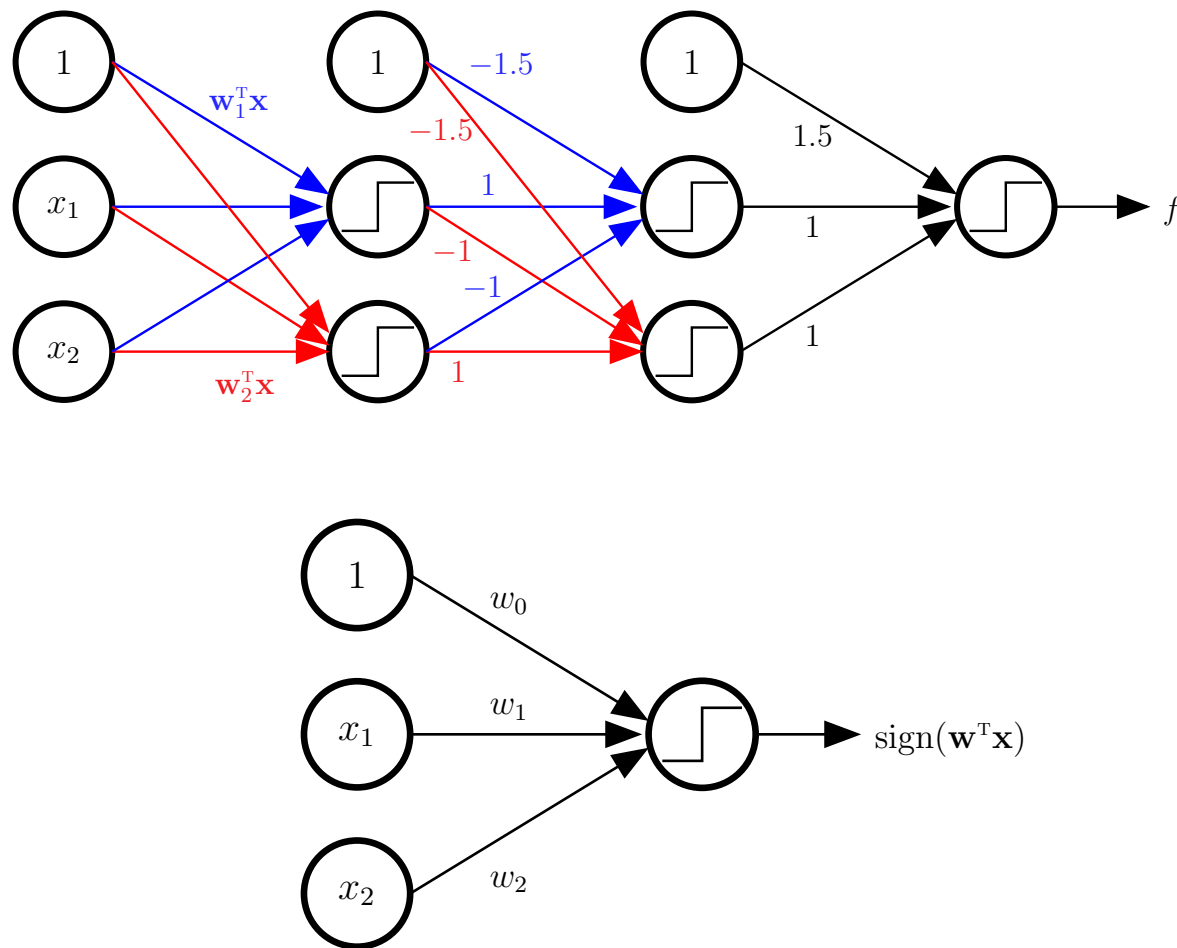


Representing f Using OR and AND

$$f = h_1 \overline{h_2} + \overline{h_1} h_2$$



The Multilayer Perceptron (MLP)



More layers allow us to implement f

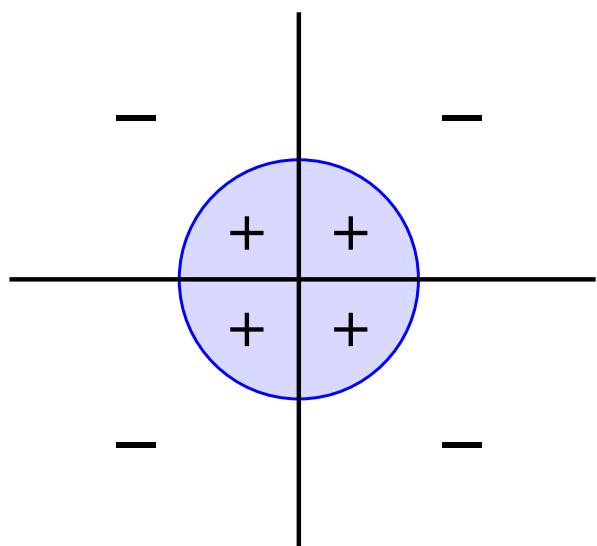
These additional layers are called *hidden layers*

Universal Approximation

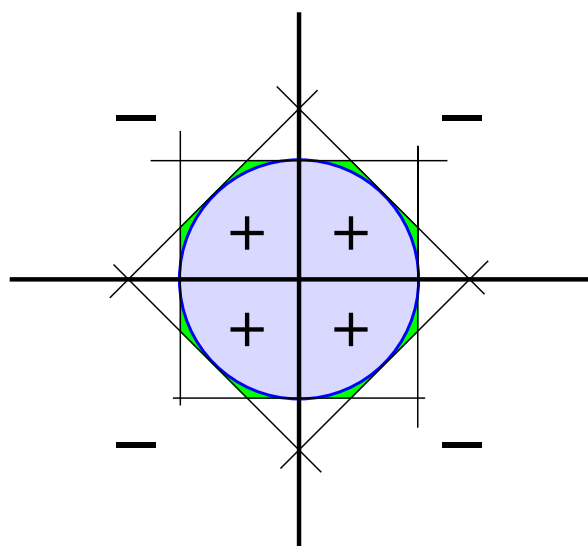
Any target function f that can be decomposed into linear separators can be implemented by a 3-layer MLP.

Universal Approximation

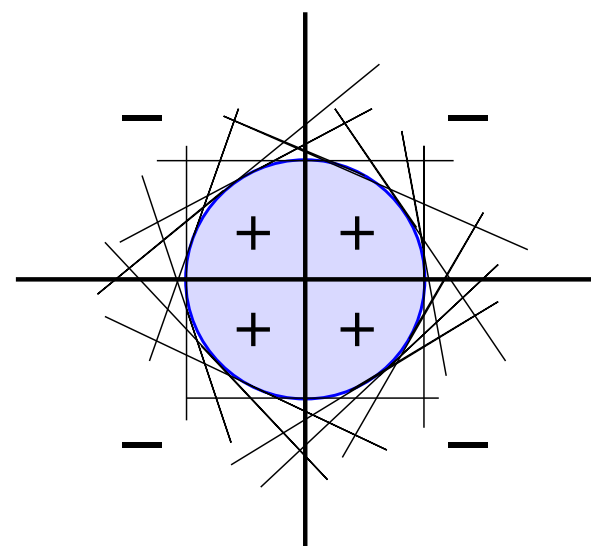
A sufficiently smooth separator can “essentially” be decomposed into linear separators.



Target



8 perceptrons



16 perceptrons

Approximation Versus Generalization

The size of the MLP controls the approximation-generalization tradeoff.

More nodes per hidden layer \implies approximation \uparrow and generalization \downarrow

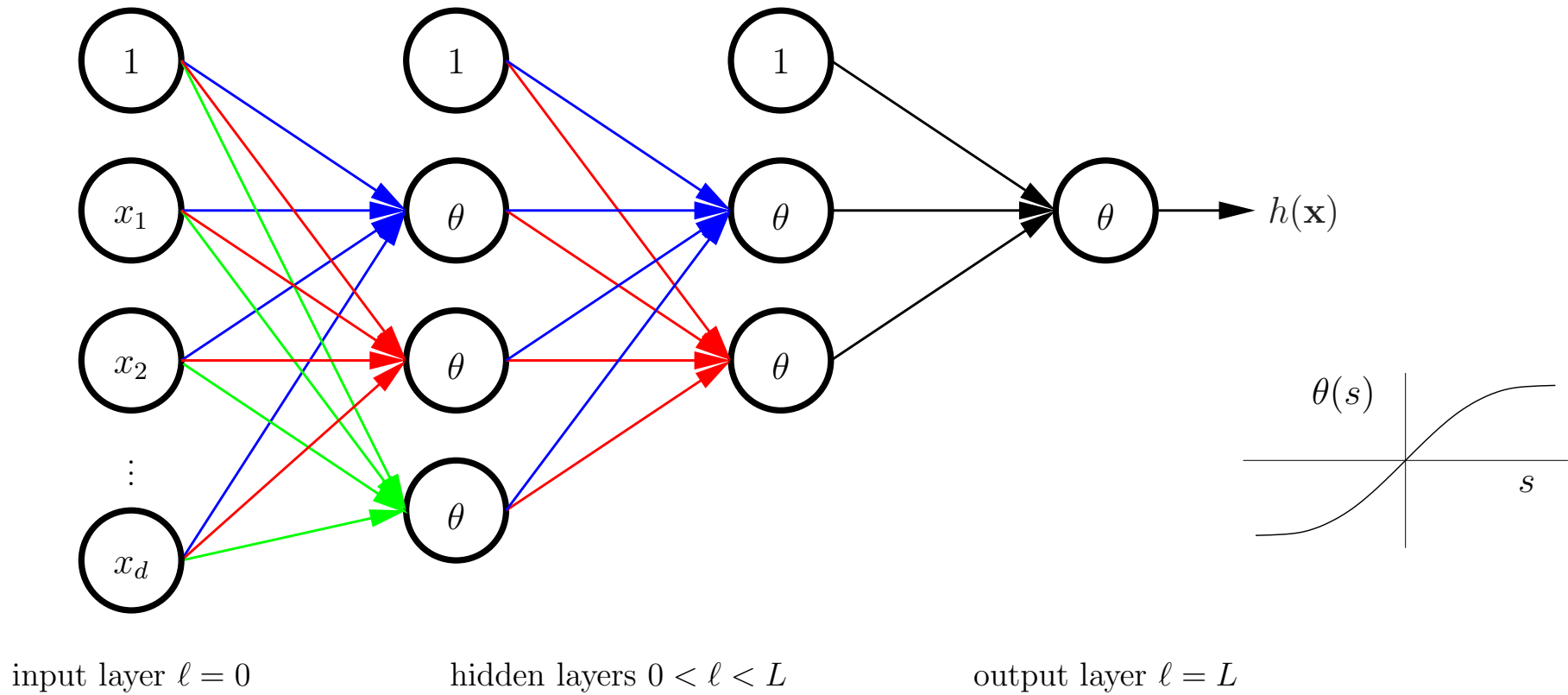
Minimizing E_{in}

A combinatorial problem even harder with the MLP than the Perceptron.

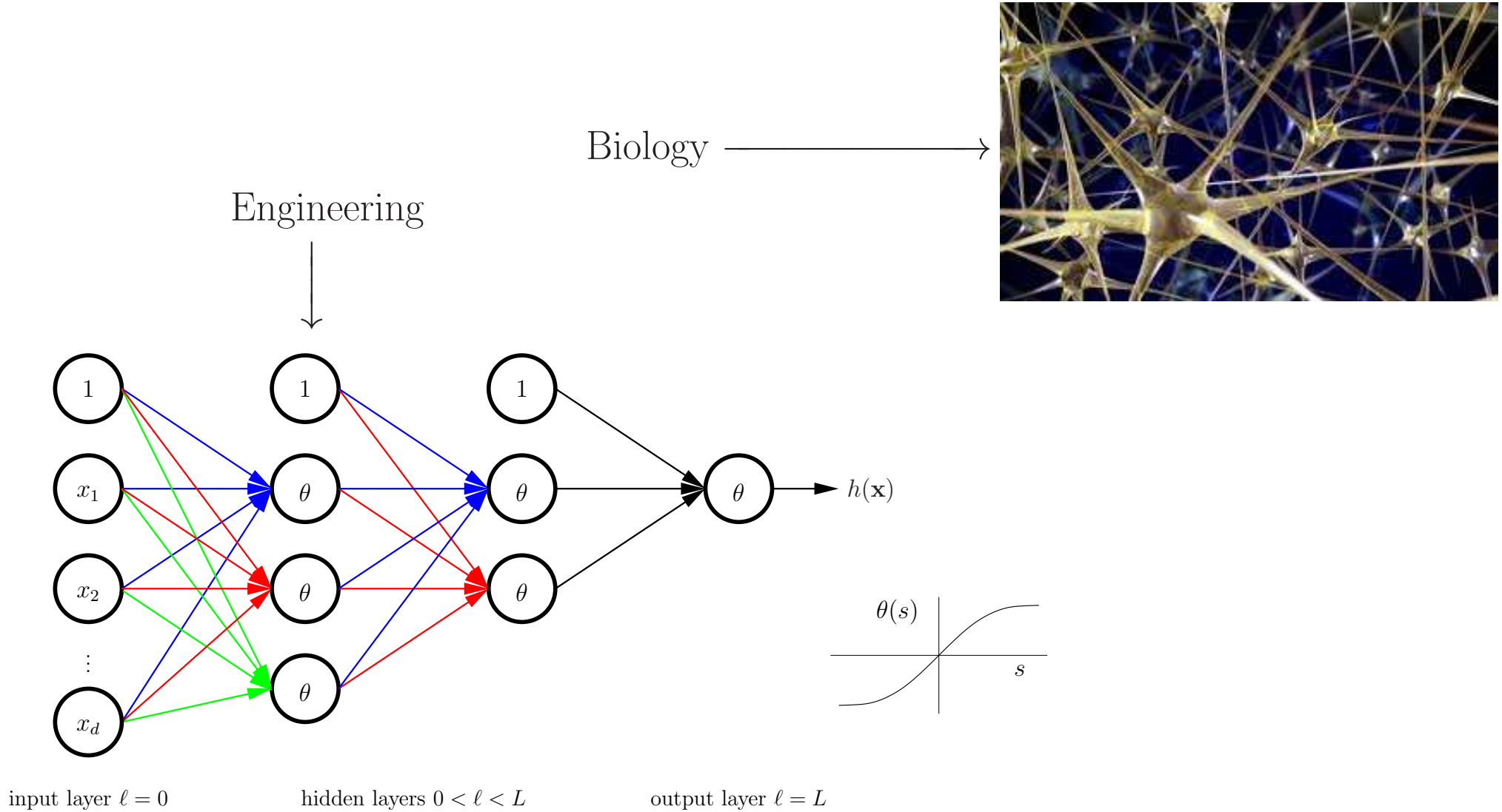
E_{in} is not smooth (due to sign function), so cannot use gradient descent.

$\text{sign}(x) \approx \tan(x) \longrightarrow$ gradient descent to minimize E_{in} .

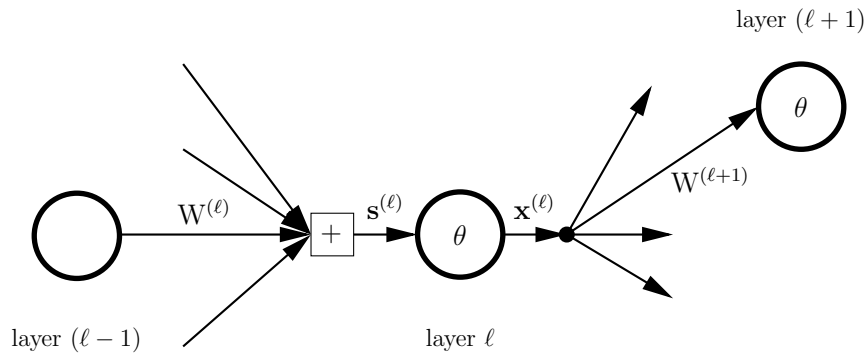
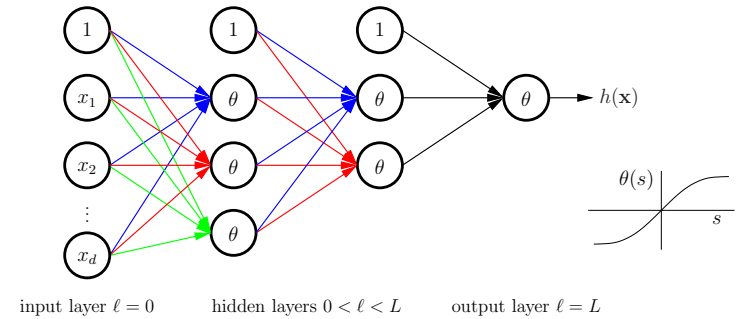
The Neural Network



The Neural Network



Zooming into a Hidden Node



layer ℓ parameters

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

layers $\ell = 0, 1, 2, \dots, L$

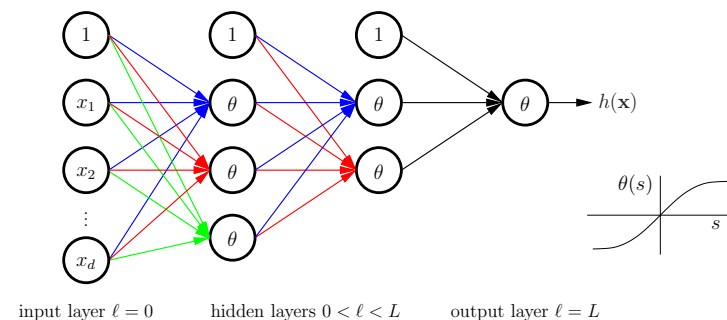
layer ℓ has “dimension” $d^{(\ell)} \implies d^{(\ell)} + 1$ nodes

$$\mathbf{W}^{(\ell)} = \begin{bmatrix} \mathbf{w}_1^{(\ell)} & \mathbf{w}_2^{(\ell)} & \cdots & \mathbf{w}_{d^{(\ell)}}^{(\ell)} \\ | & | & \vdots & | \end{bmatrix}$$

The Linear Signal

Input $\mathbf{s}^{(\ell)}$ is a linear combination (using weights) of the outputs of the previous layer $\mathbf{x}^{(\ell-1)}$.

$$\mathbf{s}^{(\ell)} = (\mathbf{W}^{(\ell)})^T \mathbf{x}^{(\ell-1)}$$



$$\begin{bmatrix} s_1^{(\ell)} \\ s_2^{(\ell)} \\ \vdots \\ s_j^{(\ell)} \\ \vdots \\ s_{d^{(\ell)}}^{(\ell)} \end{bmatrix} = \begin{bmatrix} (\mathbf{w}_1^{(\ell)})^T \text{---} \\ (\mathbf{w}_2^{(\ell)})^T \text{---} \\ \vdots \\ (\mathbf{w}_j^{(\ell)})^T \text{---} \\ \vdots \\ (\mathbf{w}_{d^{(\ell)}}^{(\ell)})^T \text{---} \end{bmatrix} \mathbf{x}^{(\ell-1)}$$

$$s_j^{(\ell)} = (\mathbf{w}_j^{(\ell)})^T \mathbf{x}^{(\ell-1)}$$

(recall the linear signal $s = \mathbf{w}^T \mathbf{x}$)

$$\mathbf{s}^{(\ell)} \xrightarrow{\theta} \mathbf{x}^{(\ell)}$$

Forward Propagation: Computing $h(\mathbf{x})$

$$\mathbf{x} = \mathbf{x}^{(0)} \xrightarrow{\mathbf{w}^{(1)}} \mathbf{s}^{(1)} \xrightarrow{\theta} \mathbf{x}^{(1)} \xrightarrow{\mathbf{w}^{(2)}} \mathbf{s}^{(2)} \xrightarrow{\theta} \mathbf{x}^{(2)} \dots \xrightarrow{\mathbf{w}^{(L)}} \mathbf{s}^{(L)} \xrightarrow{\theta} \mathbf{x}^{(L)} = h(\mathbf{x}).$$

Forward propagation to compute $h(\mathbf{x})$:

1: $\mathbf{x}^{(0)} \leftarrow \mathbf{x}$

[Initialization]

2: **for** $\ell = 1$ to L **do**

[Forward Propagation]

3: $\mathbf{s}^{(\ell)} \leftarrow (\mathbf{W}^{(\ell)})^T \mathbf{x}^{(\ell-1)}$

4: $\mathbf{x}^{(\ell)} \leftarrow \begin{bmatrix} 1 \\ \theta(\mathbf{s}^{(\ell)}) \end{bmatrix}$

5: **end for**

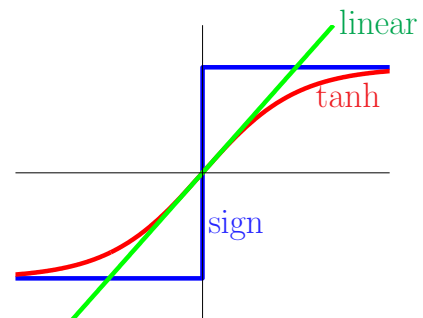
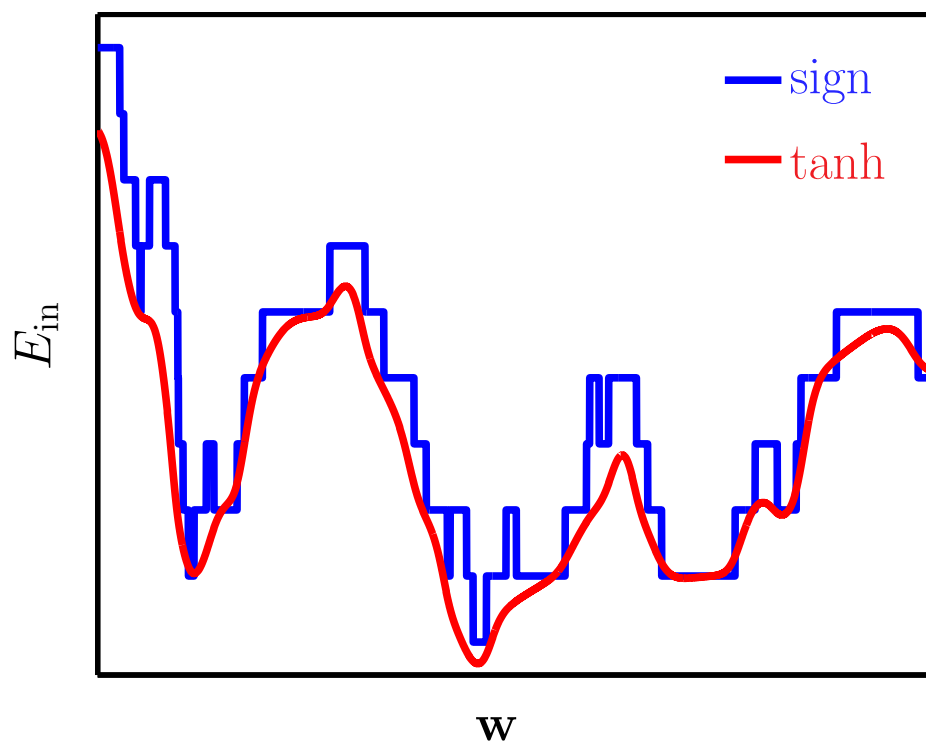
6: $h(\mathbf{x}) = \mathbf{x}^{(L)}$

[Output]

Minimizing E_{in}

$$E_{\text{in}}(h) = E_{\text{in}}(\mathbf{W}) = \frac{1}{N} \sum_{n=1}^N (h(\mathbf{x}_n) - y_n)^2$$

$$\mathbf{W} = \{\mathbf{W}^{(1)}, \mathbf{W}^{(2)}, \dots, \mathbf{W}^{(L)}\}$$




Using $\theta = \text{tanh}$ makes E_{in} differentiable so we can use gradient descent \rightarrow local minimum.

Gradient Descent

$$W(t + 1) = W(t) - \eta \nabla E_{\text{in}}(W(t))$$

Gradient of E_{in}

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \mathbf{e}(h(\mathbf{x}_n), y_n)$$


$$\frac{\partial E_{\text{in}}(\mathbf{w})}{\partial W^{(\ell)}} = \frac{1}{N} \sum_{n=1}^N \frac{\partial \mathbf{e}_n}{\partial W^{(\ell)}}$$

We need

$$\frac{\partial \mathbf{e}(\mathbf{x})}{\partial W^{(\ell)}}$$

Numerical Approach

$$\frac{\partial \mathbf{e}(\mathbf{x})}{\partial W_{ij}^{(\ell)}} \approx \frac{\mathbf{e}(\mathbf{x}|W_{ij}^{(\ell)} + \Delta) - \mathbf{e}(\mathbf{x}|W_{ij}^{(\ell)} - \Delta)}{2\Delta}$$

approximate

inefficient



Example Data

