



# CS 540 Introduction to Artificial Intelligence Neural Networks (II)

Yudong Chen

University of Wisconsin-Madison

Oct 21, 2021

# Announcement

- Midterm: Oct 28 (Thursday)
  - Details can be found on [https://piazza.com/class/ktahk2n2c6wkp?  
cid=161](https://piazza.com/class/ktahk2n2c6wkp?cid=161)

# Today's outline

- Single-layer Perceptron: Recap
- Multi-layer Perceptron
  - Single output
  - Multiple output
- How to train neural networks
  - Gradient descent

# Review: Perceptron

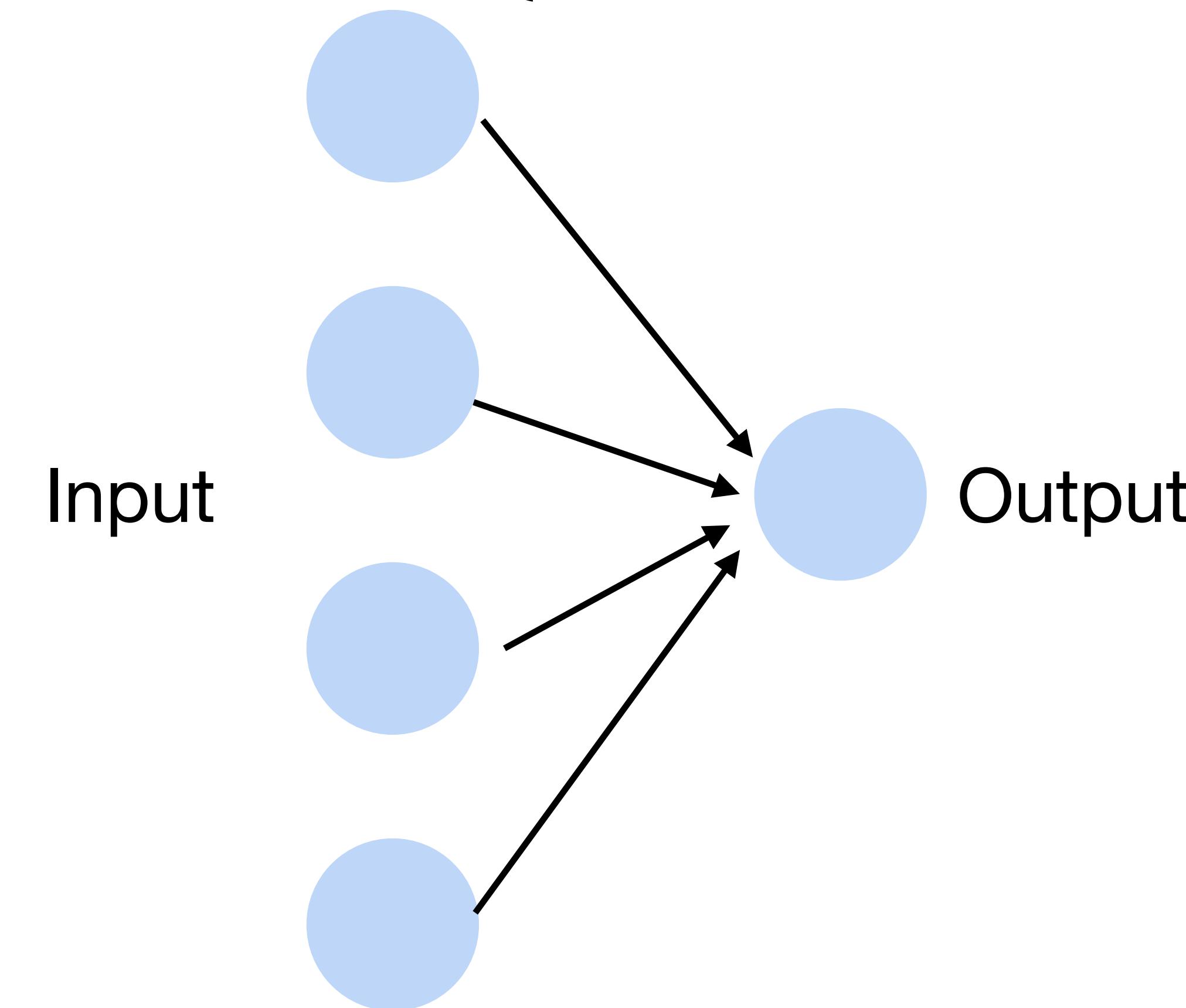
- Given input  $\mathbf{x}$ , weight  $\mathbf{w}$  and bias  $b$ , perceptron outputs:

$$f(\mathbf{x}) = \sigma(\langle \mathbf{w}, \mathbf{x} \rangle + b)$$

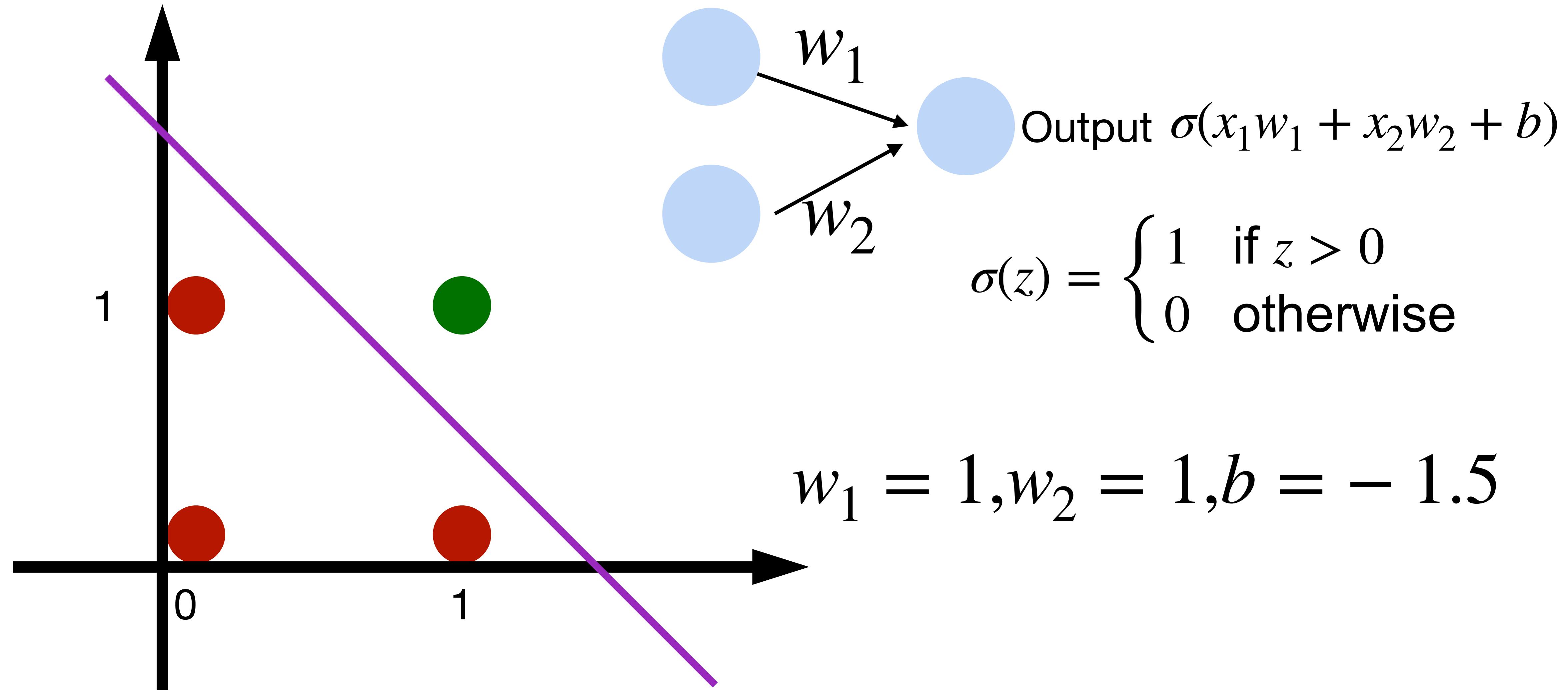
$$\sigma(z) = \begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{otherwise} \end{cases}$$

Activation function

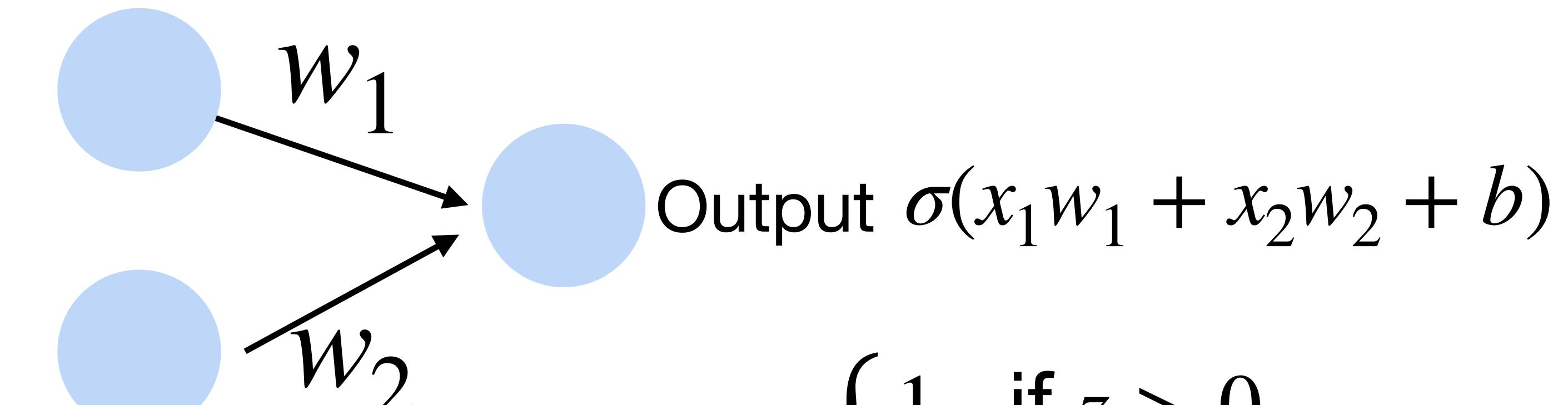
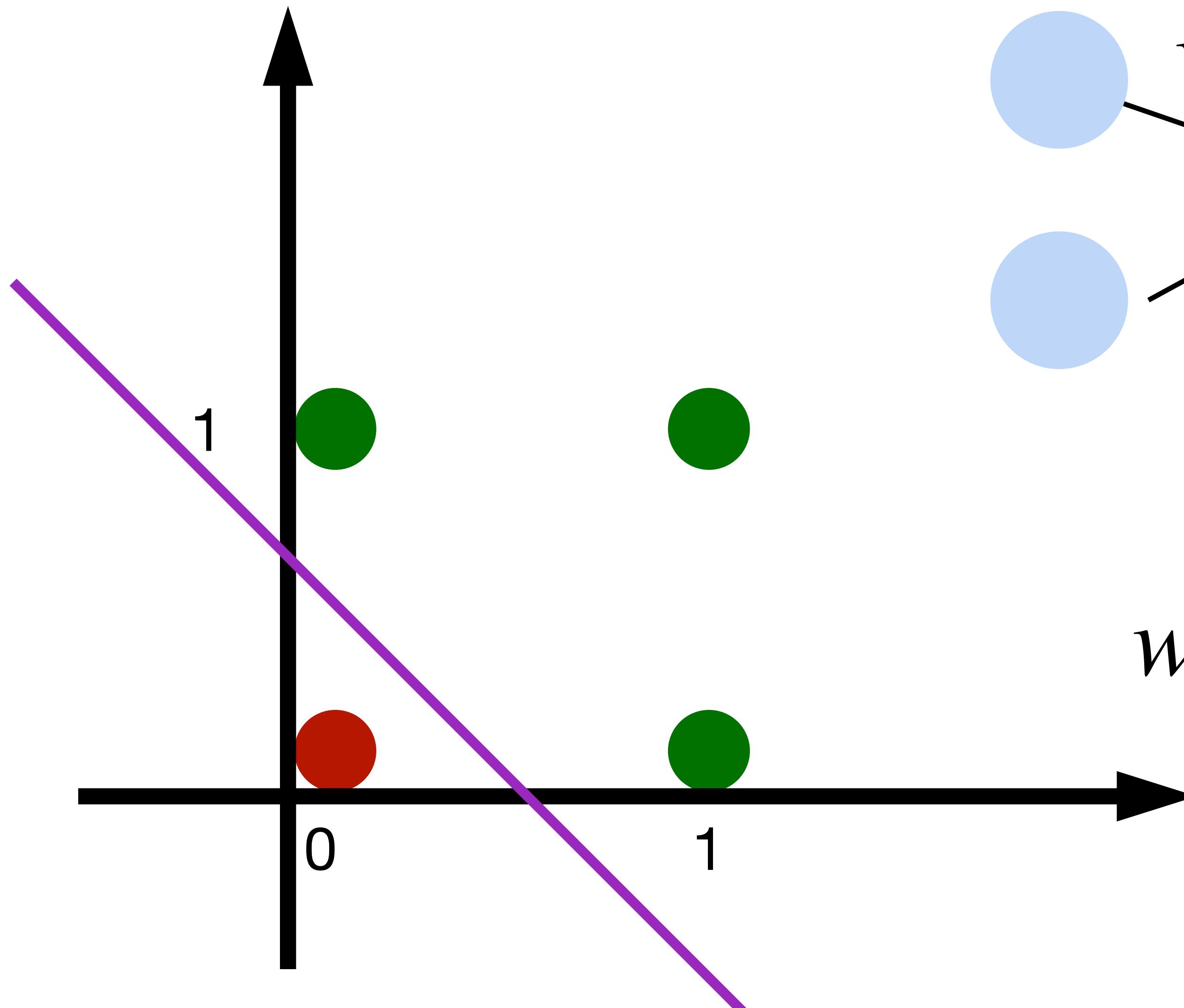
Cats vs. dogs?



# Perceptron can learn AND function



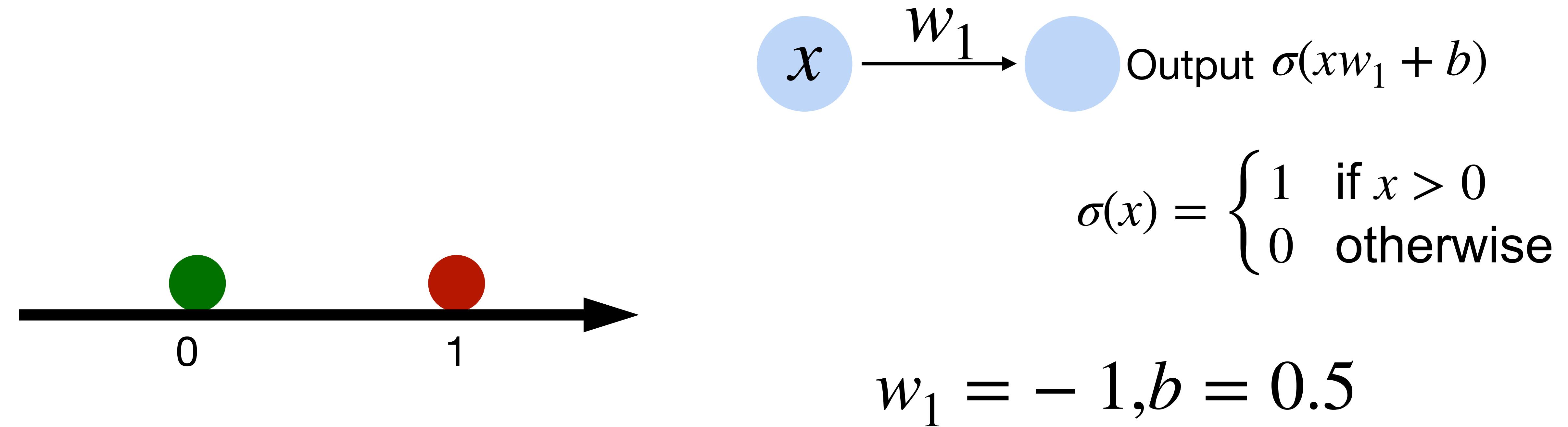
# Perceptron can learn OR function



$$\sigma(z) = \begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{otherwise} \end{cases}$$

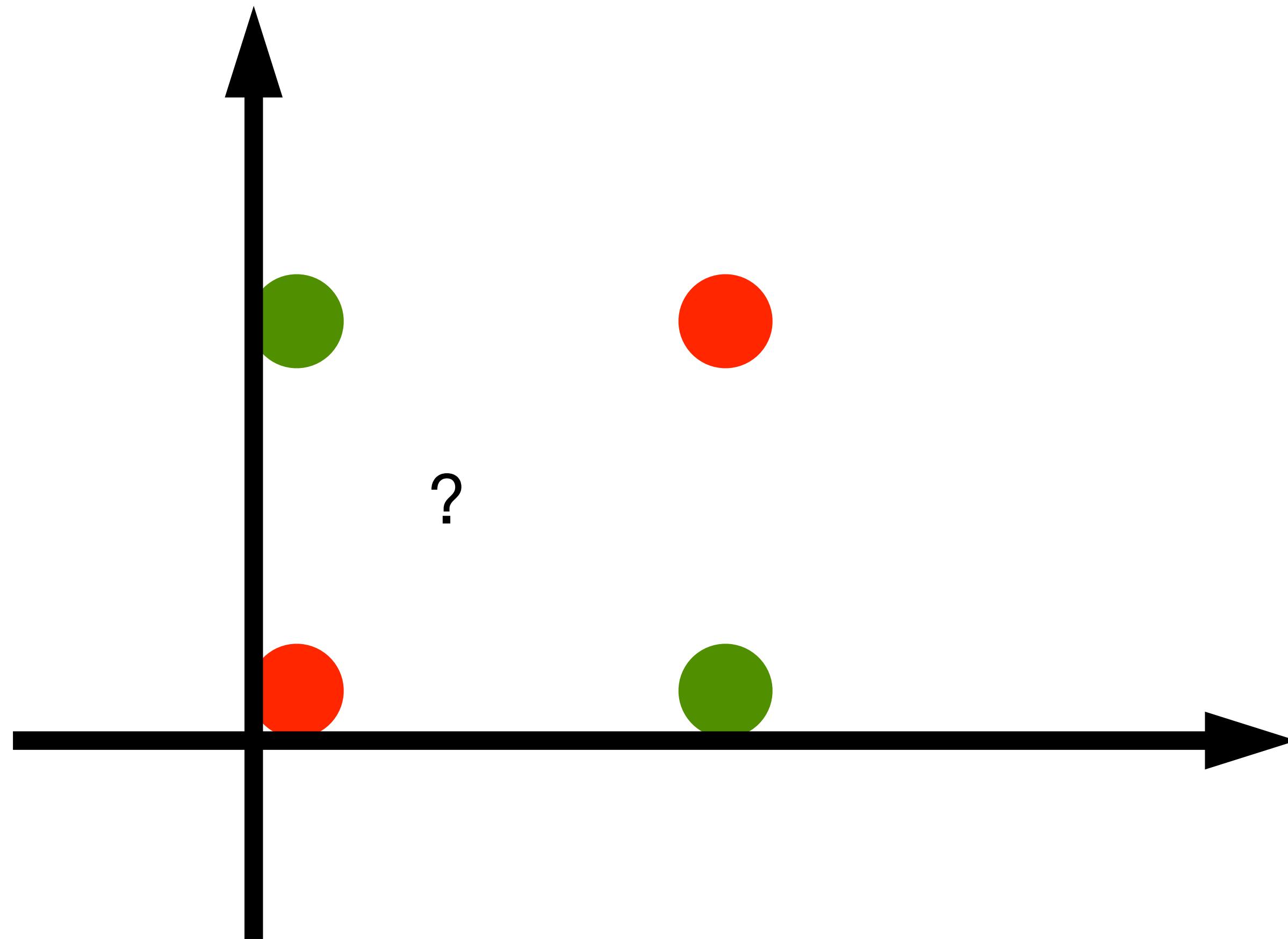
$$w_1 = 1, w_2 = 1, b = -0.5$$

# Perceptron can learn NOT function



# The limited power of a single neuron

The perceptron cannot learn an **XOR** function  
(neurons can only generate linear separators)



$$x_1 = 1, x_2 = 1, y = 0$$

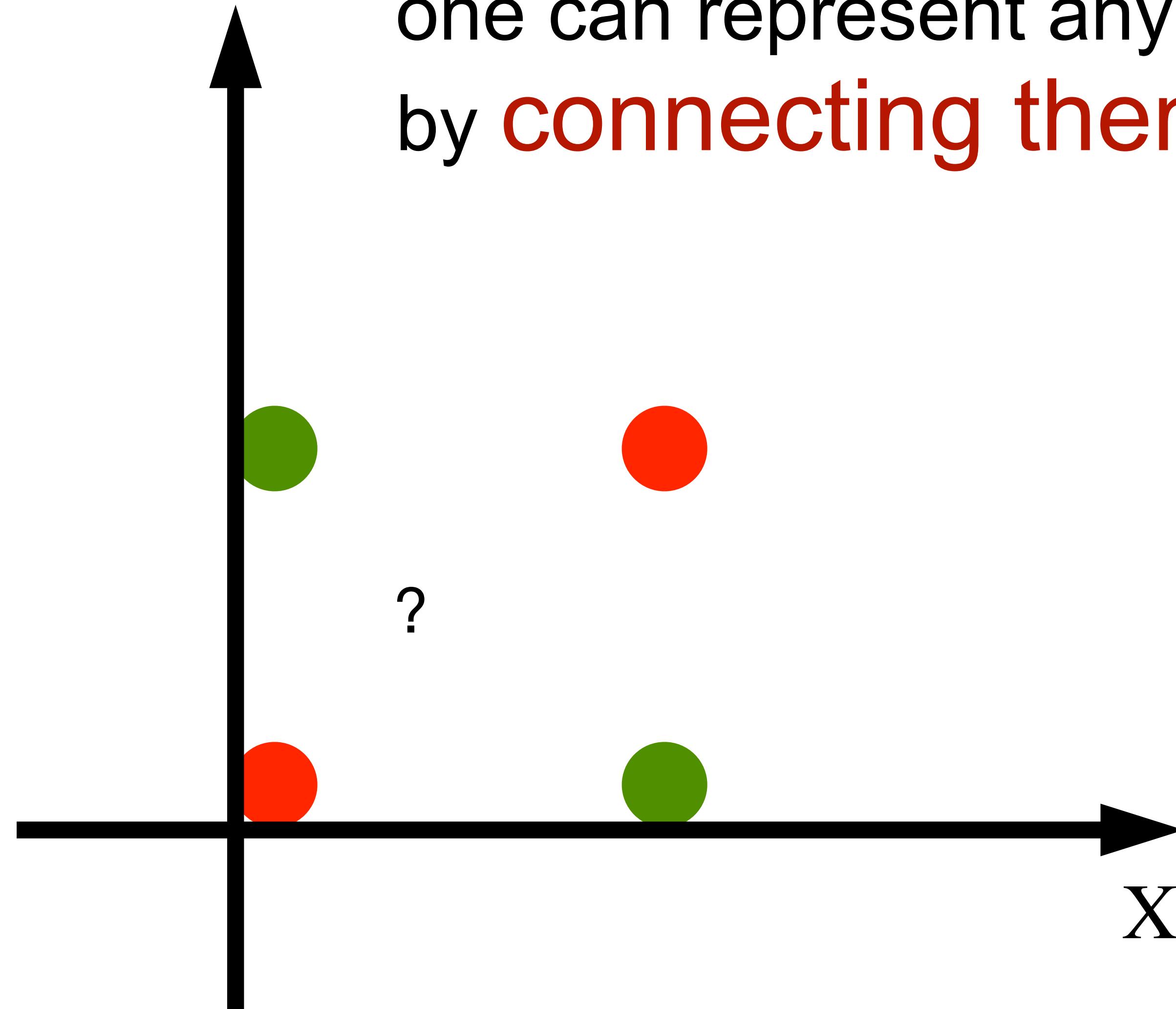
$$x_1 = 1, x_2 = 0, y = 1$$

$$x_1 = 0, x_2 = 1, y = 1$$

$$x_1 = 0, x_2 = 0, y = 0$$

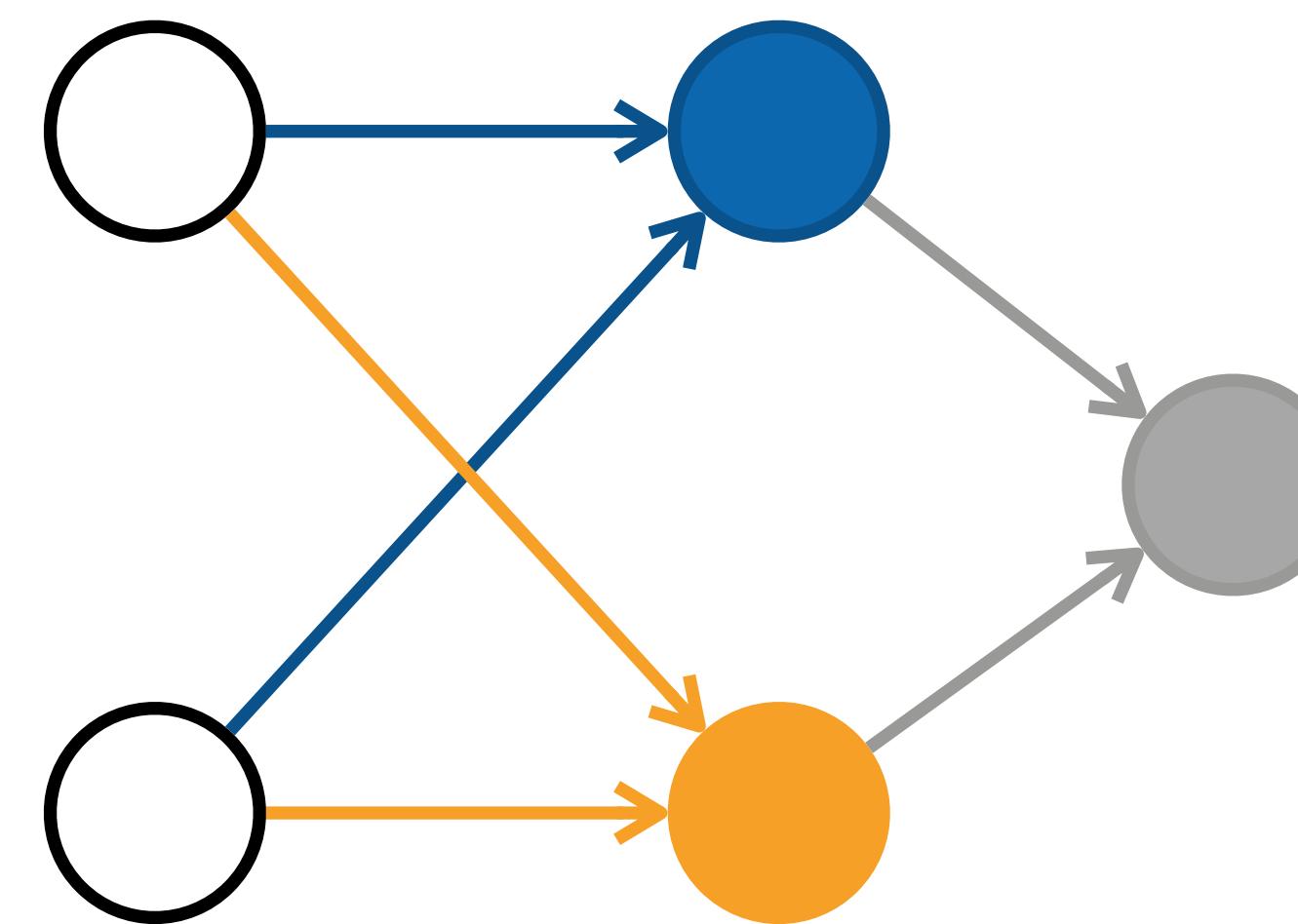
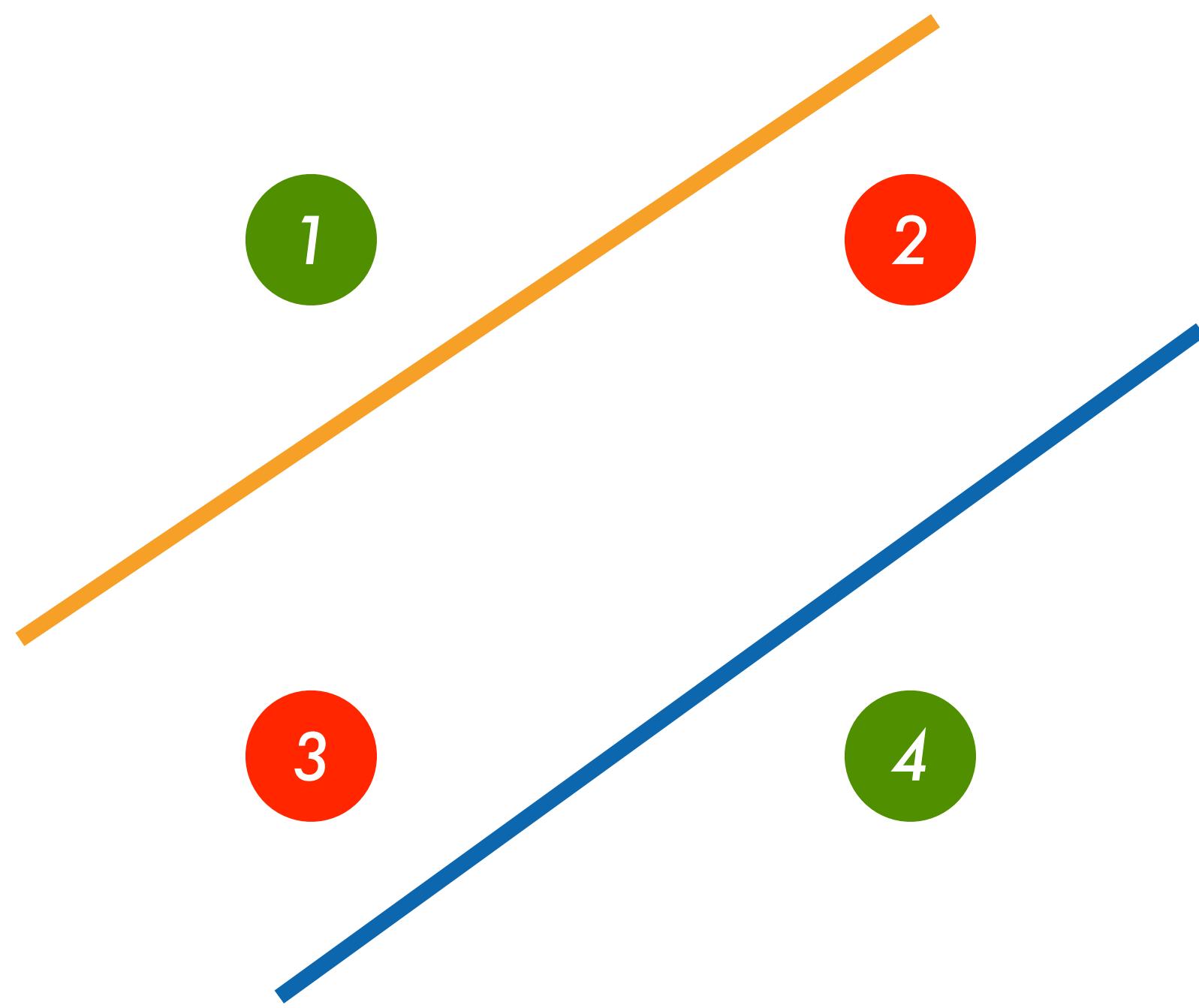
# Overcoming the limit of a single neuron

If one can represent AND, OR, NOT,  
one can represent any logic circuit (including XOR),  
by **connecting them**

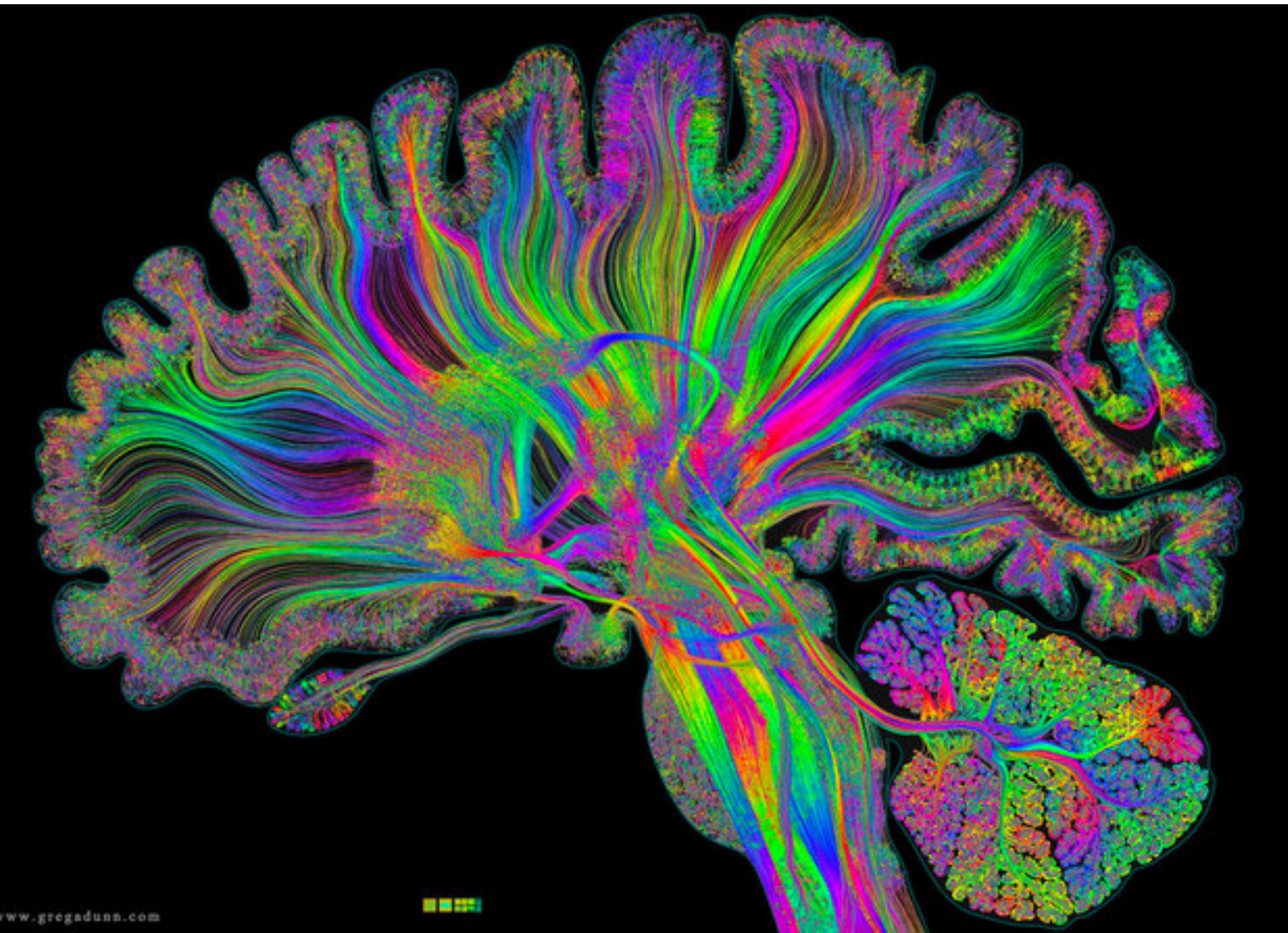


$$\text{XOR}(x_1, x_2) = (x_1 \wedge \neg x_2) \vee (\neg x_1 \wedge x_2)$$

# Learning XOR

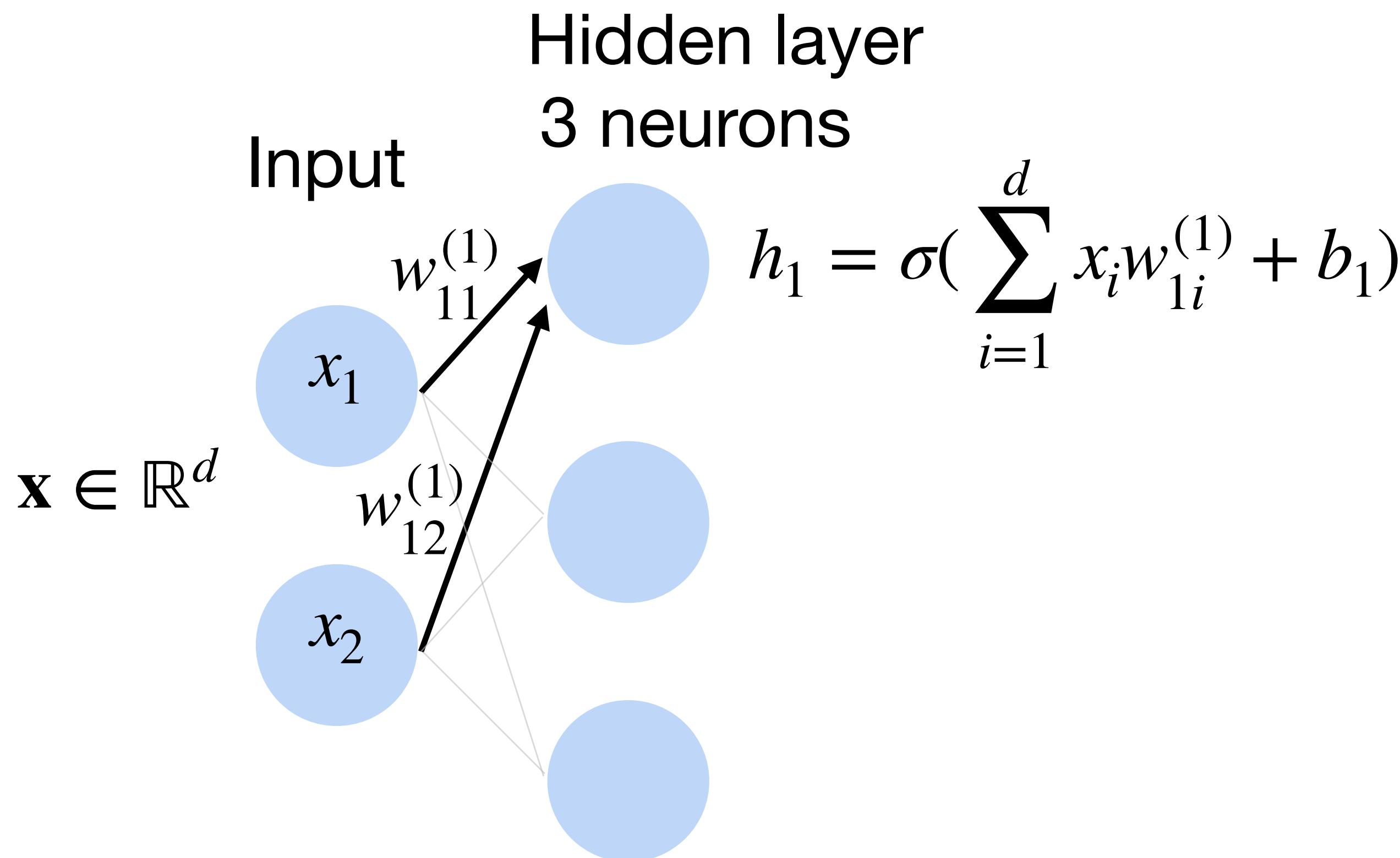


# Multilayer Perceptron



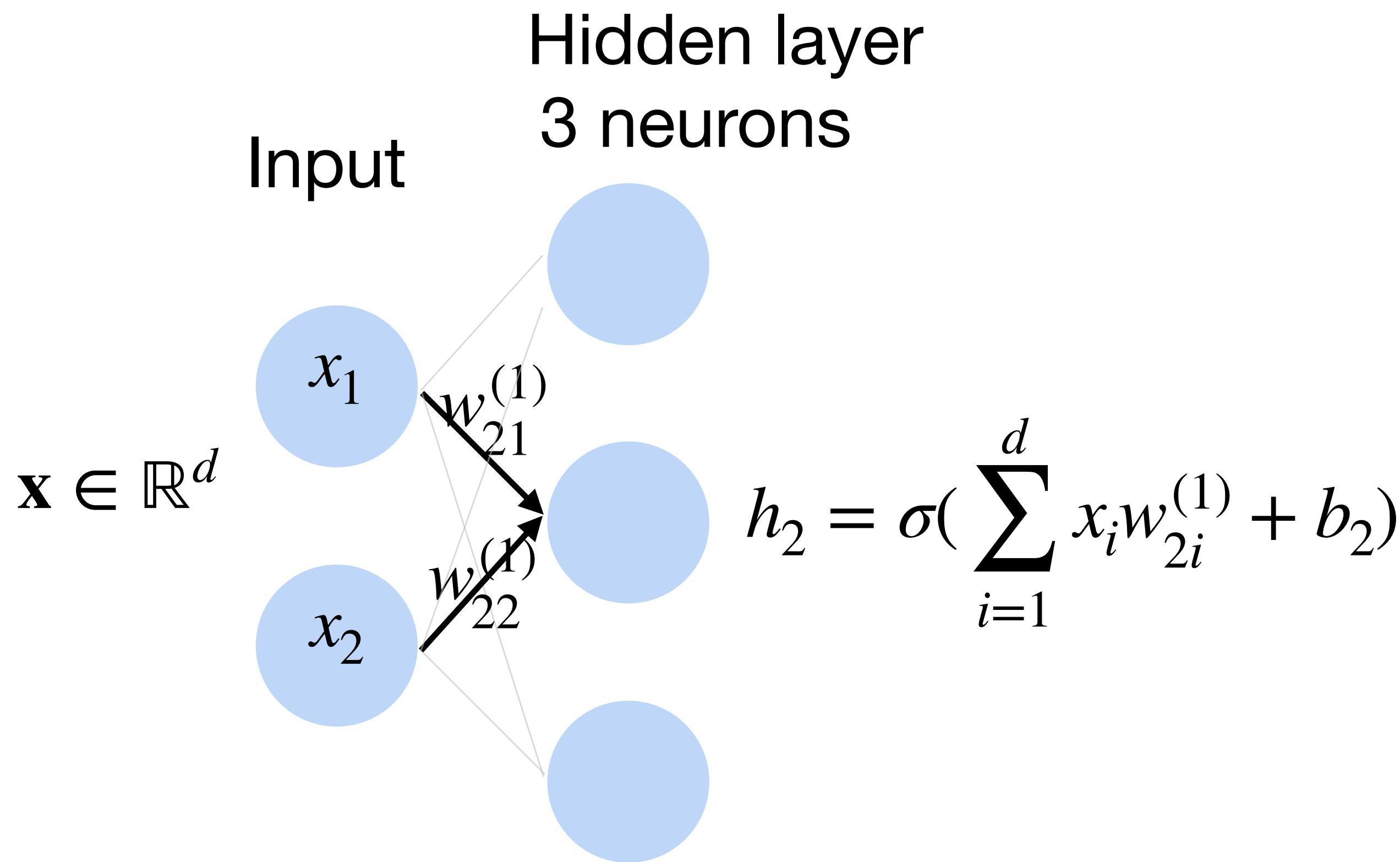
# Multi-layer perceptron: Example

- Example: 1 hidden layer, 1 output layer (depth = 2)



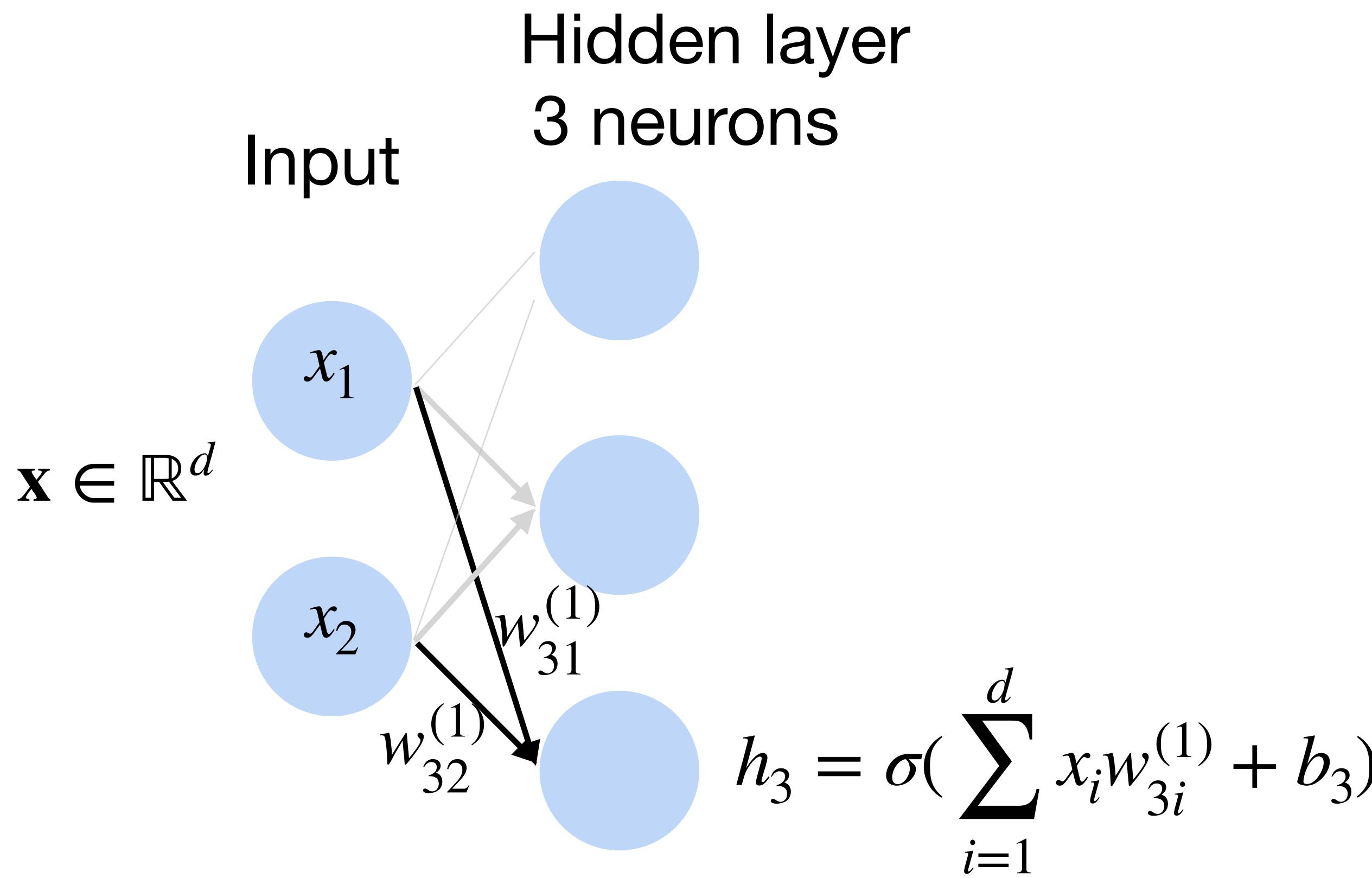
# Multi-layer perceptron: Example

- Example: 1 hidden layer, 1 output layer (depth = 2)



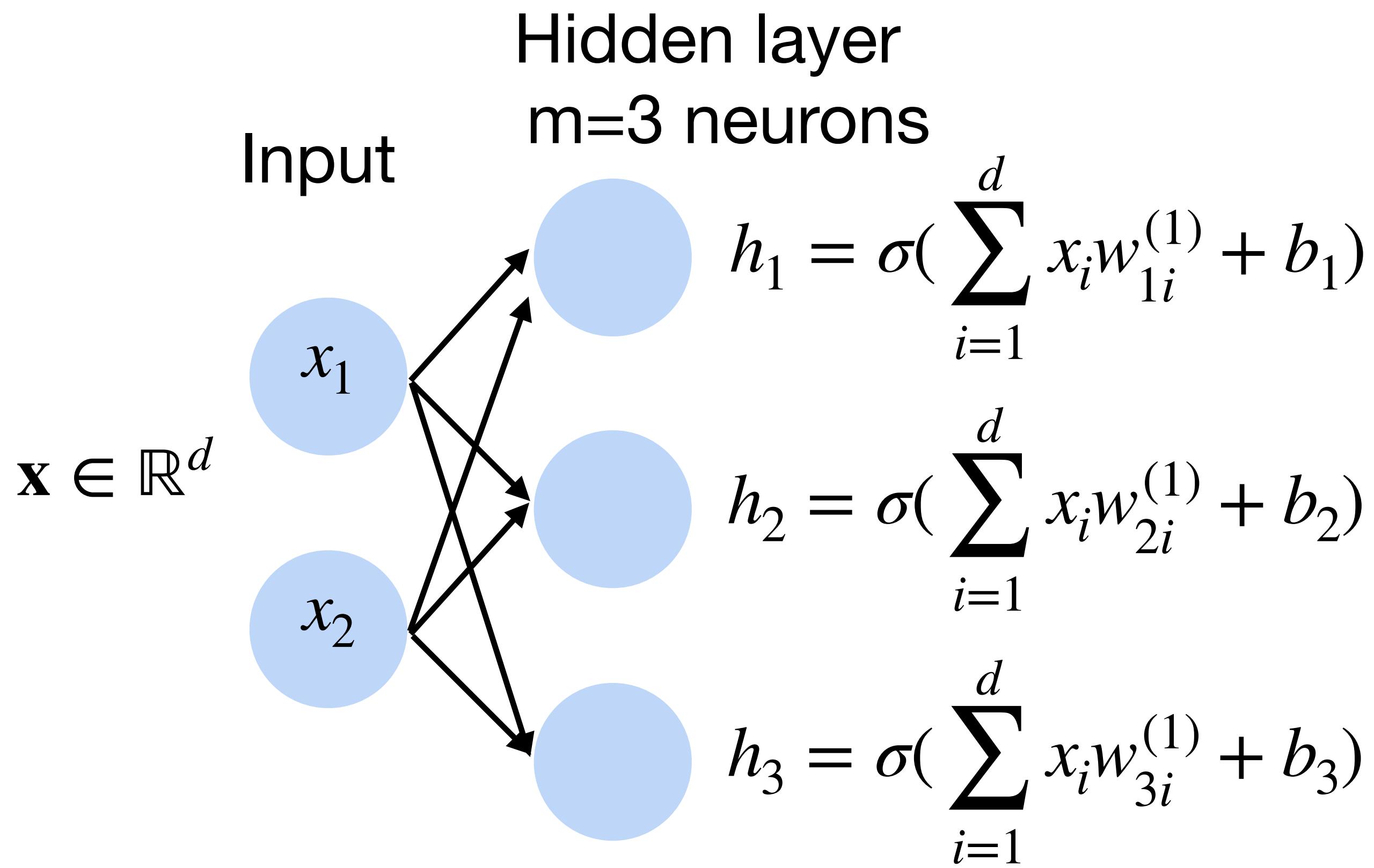
# Multi-layer perceptron: Example

- Example: 1 hidden layer, 1 output layer (depth = 2)



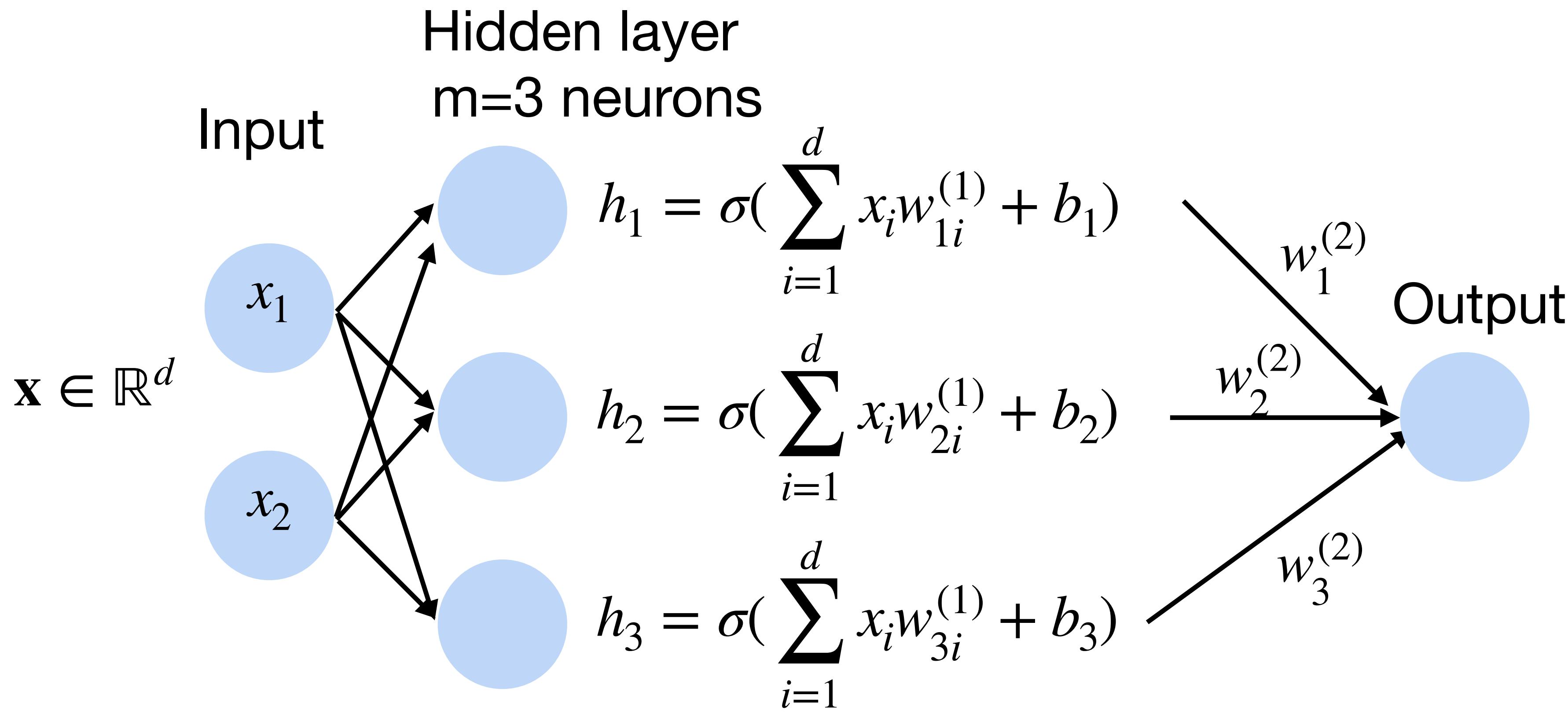
# Multi-layer perceptron: Example

- Example: 1 hidden layer, 1 output layer (depth = 2)



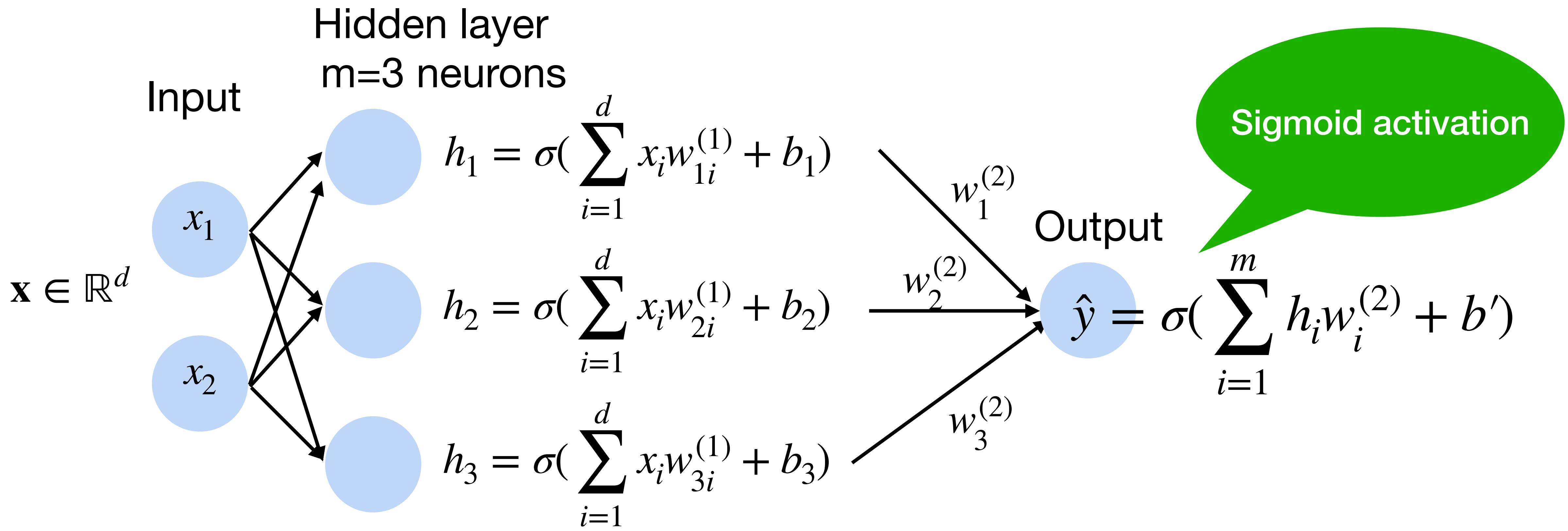
# Multi-layer perceptron: Example

- Example: 1 hidden layer, 1 output layer (depth = 2)



# Multi-layer perceptron: Example

- Example: 1 hidden layer, 1 output layer (depth = 2)



# Multi-layer perceptron: Matrix Notation

- Input  $\mathbf{x} \in \mathbb{R}^d$
- Parameters for hidden layer

$$\mathbf{W}^{(1)} \in \mathbb{R}^{m \times d}, \mathbf{b} \in \mathbb{R}^m$$

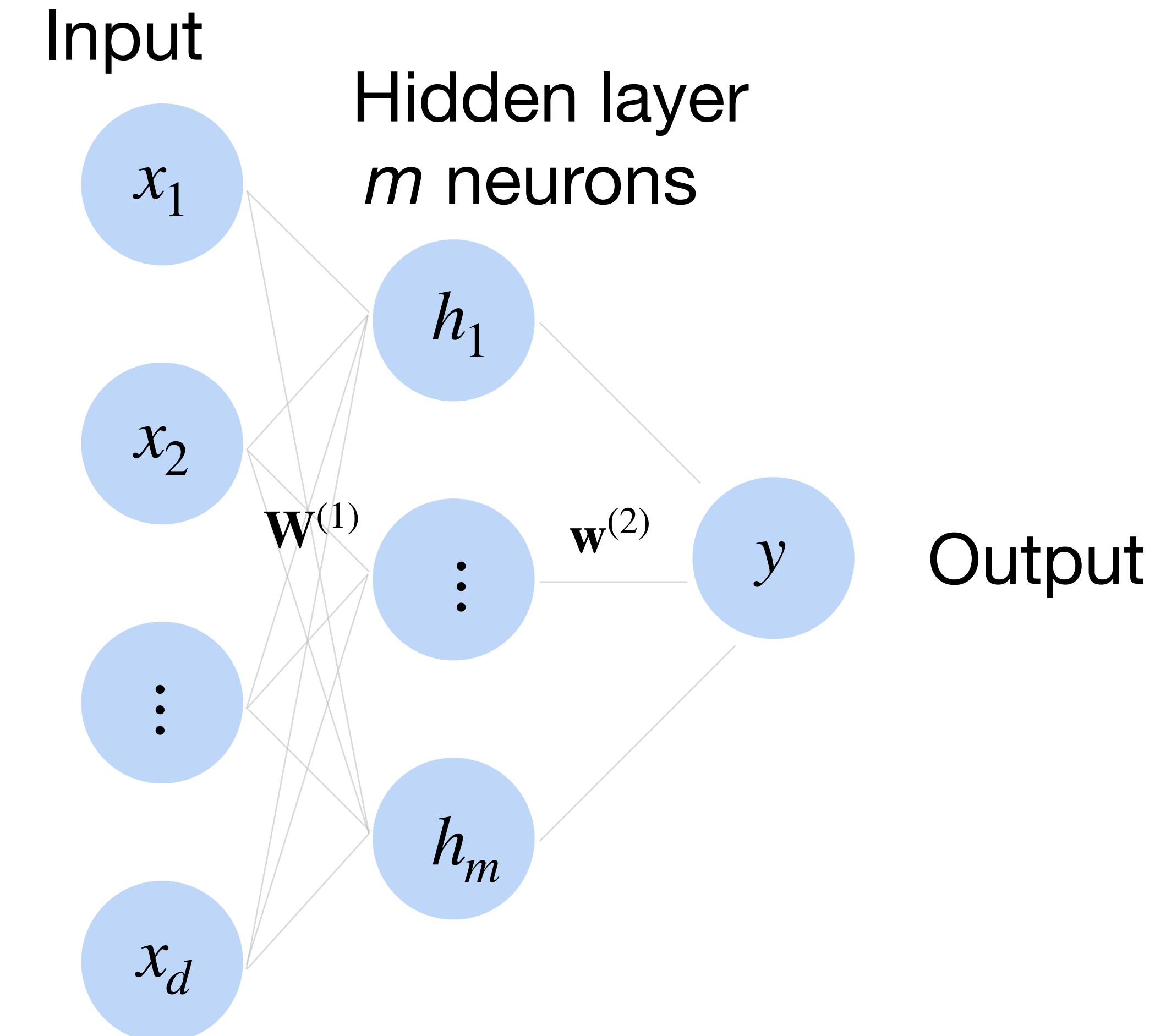
- Output from hidden layer

$$\mathbf{h} = \sigma(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b})$$

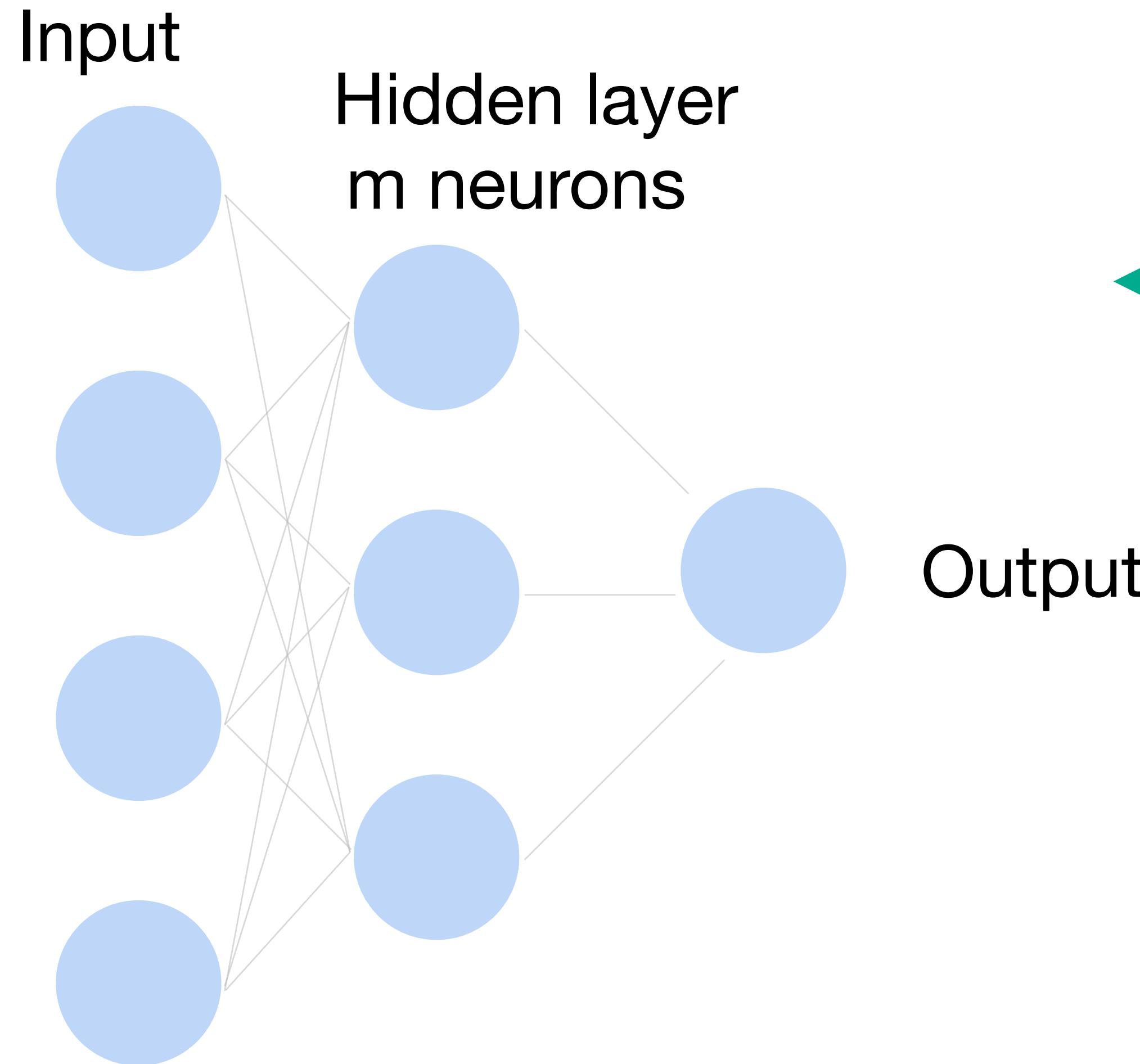
$$\mathbf{h} \in \mathbb{R}^m$$

- Final output

$$y = \mathbf{w}^{(2)T}\mathbf{h} + b'$$

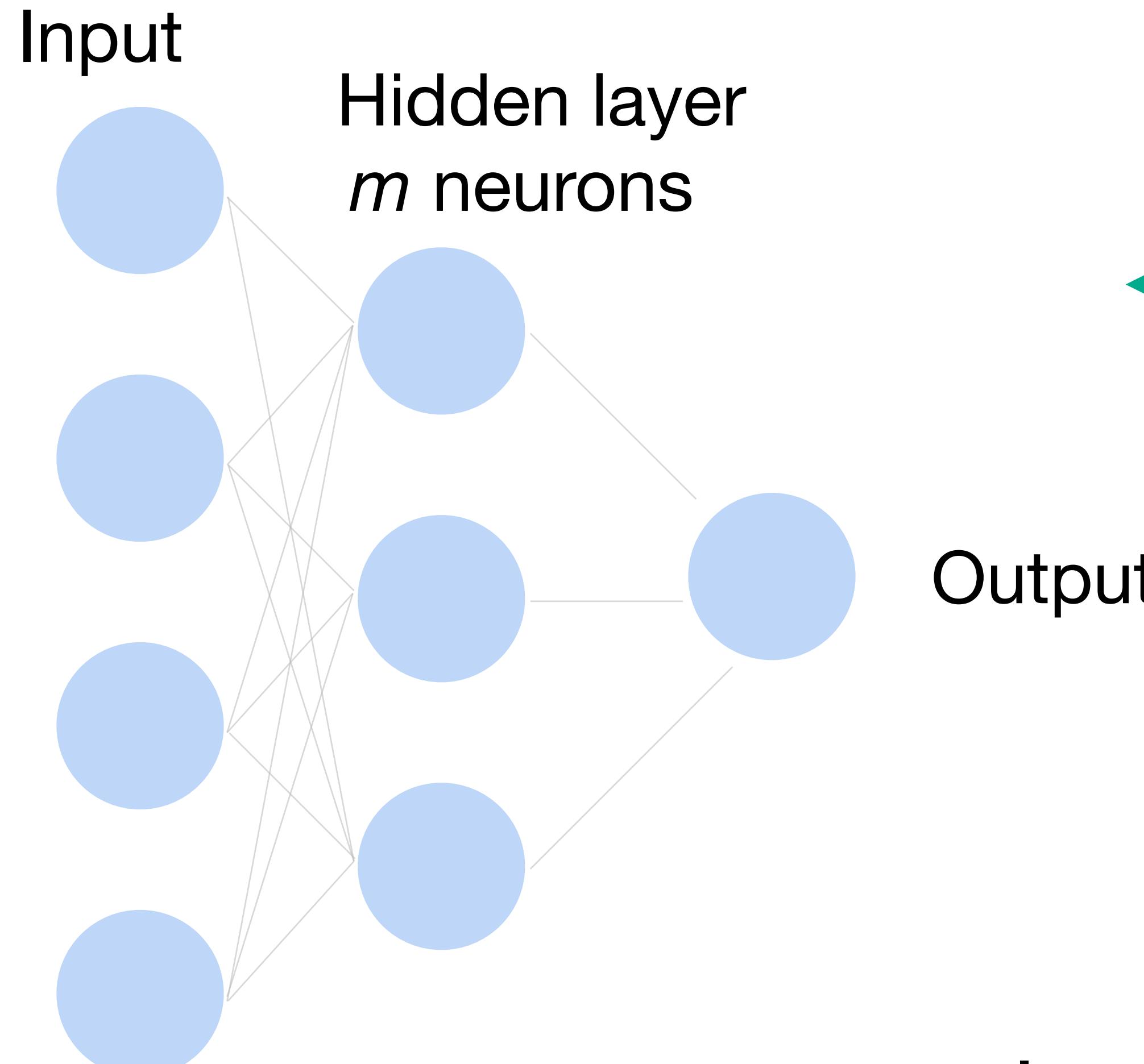


# Multi-layer perceptron



Why do we need an a  
nonlinear activation?

# Multi-layer perceptron



Why do we need an a  
nonlinear activation?

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

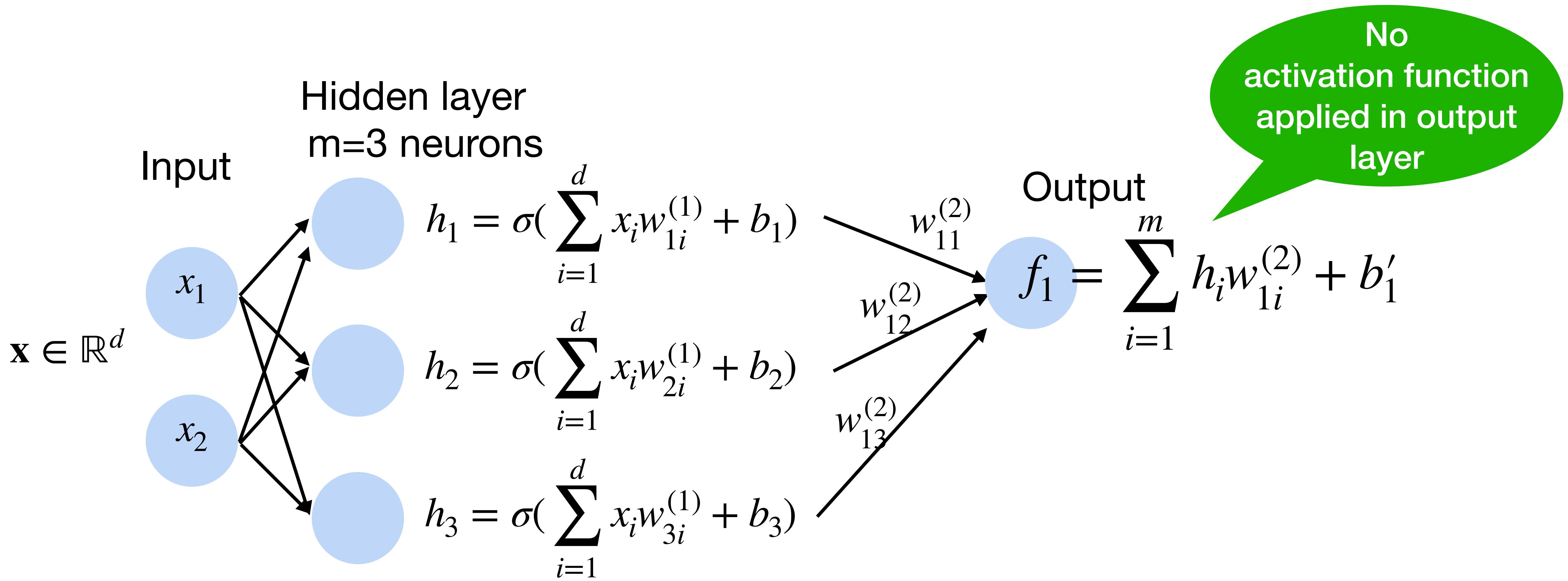
$$f(\mathbf{x}) = \mathbf{w}^{(2)\top}\mathbf{h} + b'$$

•

$$\text{hence } f(\mathbf{x}) = \mathbf{w}^{(2)\top}\mathbf{W}^{(1)}\mathbf{x} + b'$$

# Neural network for $k$ -way classification

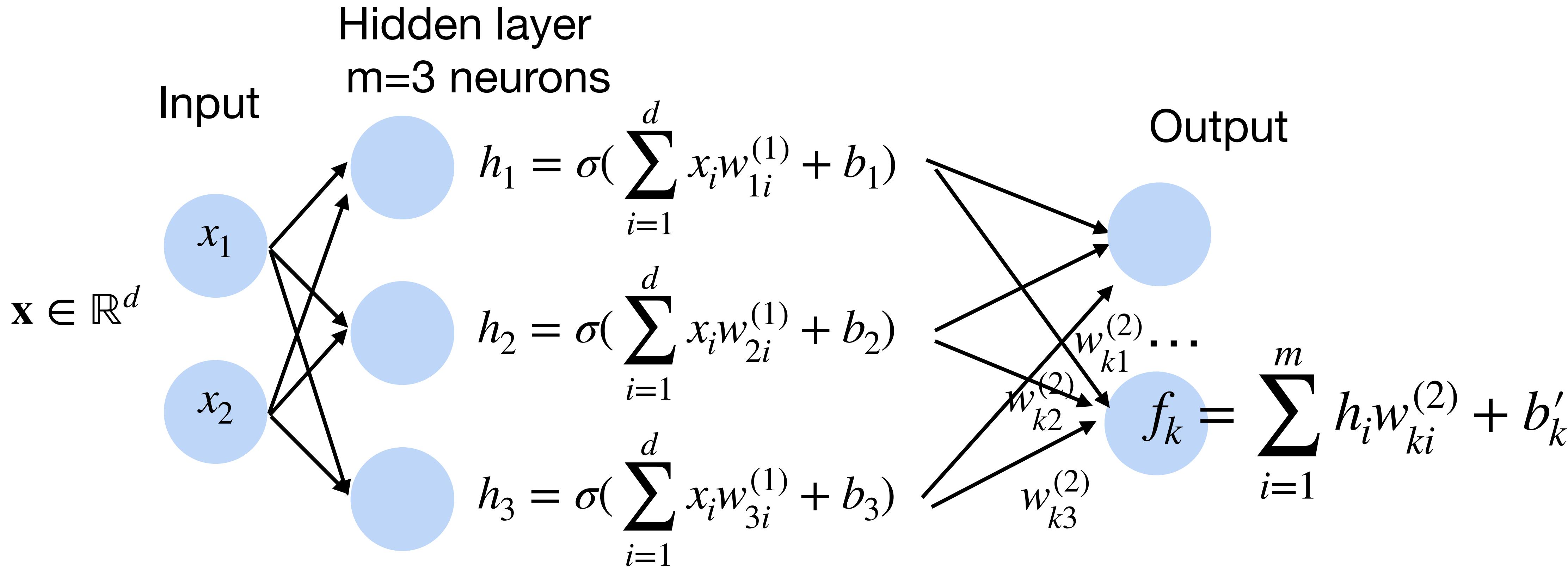
- $k$  outputs in the final layer



# Neural network for $k$ -way classification

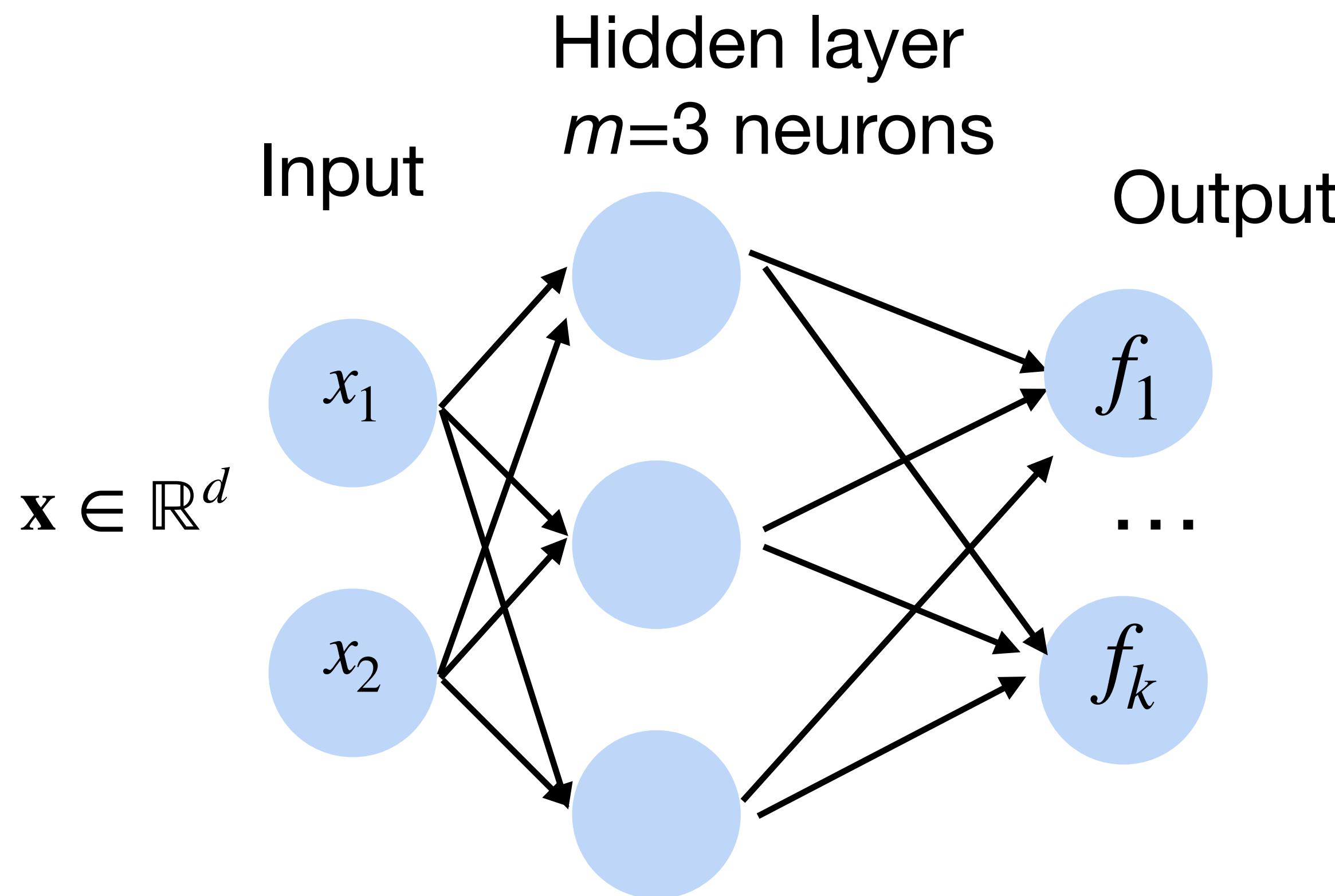
- $k$  outputs units in the final layer

**Multi-class classification** (e.g., ImageNet with  $k=1000$ )



# Softmax

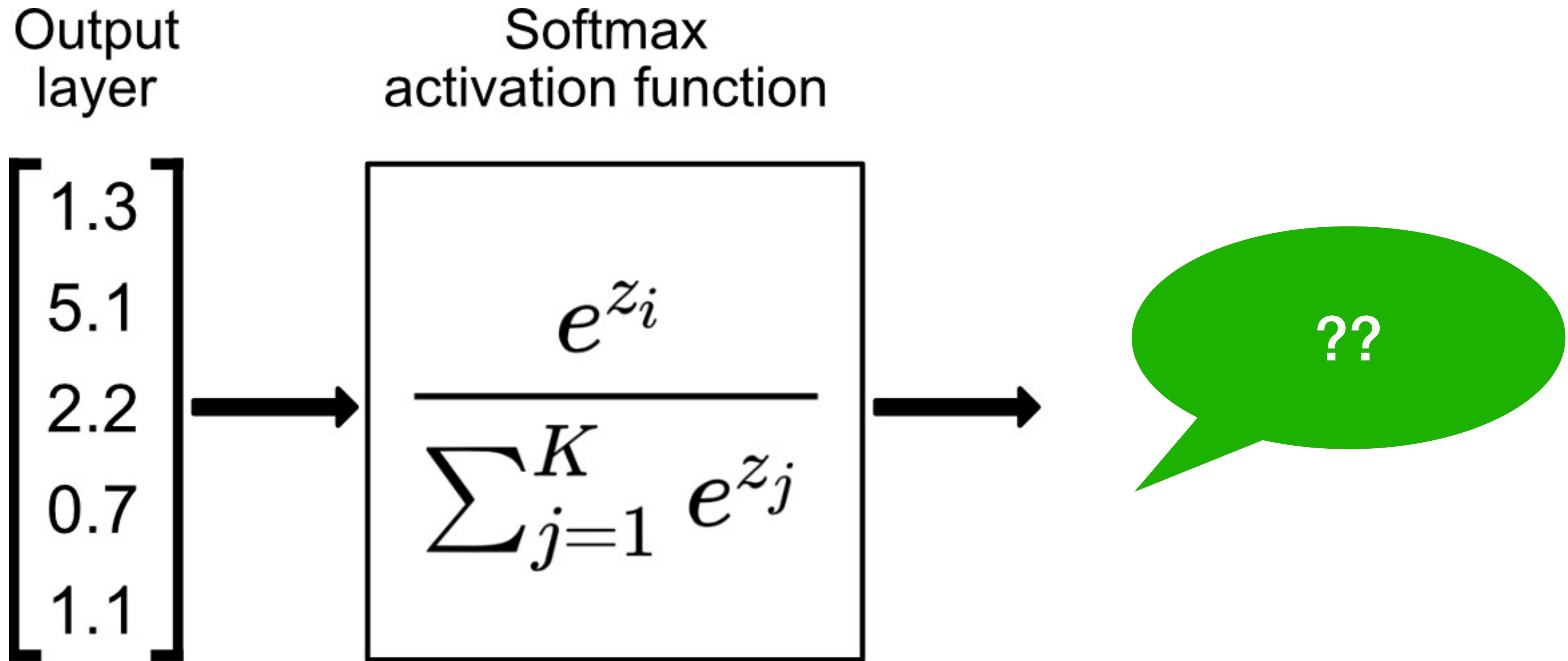
Turns outputs  $f$  into probabilities (sum up to 1 across  $k$  classes)



$$p(y = i | \mathbf{x}) = \text{softmax}(\vec{f})$$
$$= \frac{\exp[f_i(x)]}{\sum_{j=1}^k \exp[f_j(x)]}$$

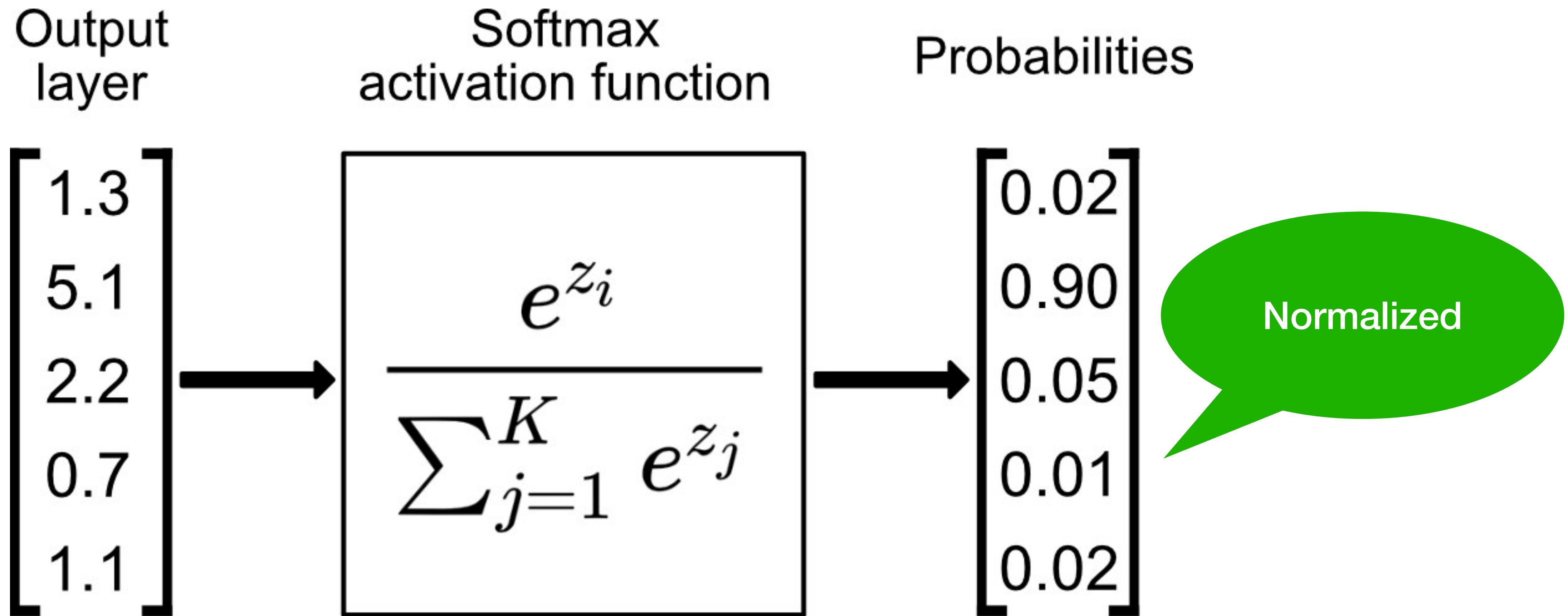
# Softmax

Turns outputs  $f$  into probabilities (sum up to 1 across  $k$  classes)



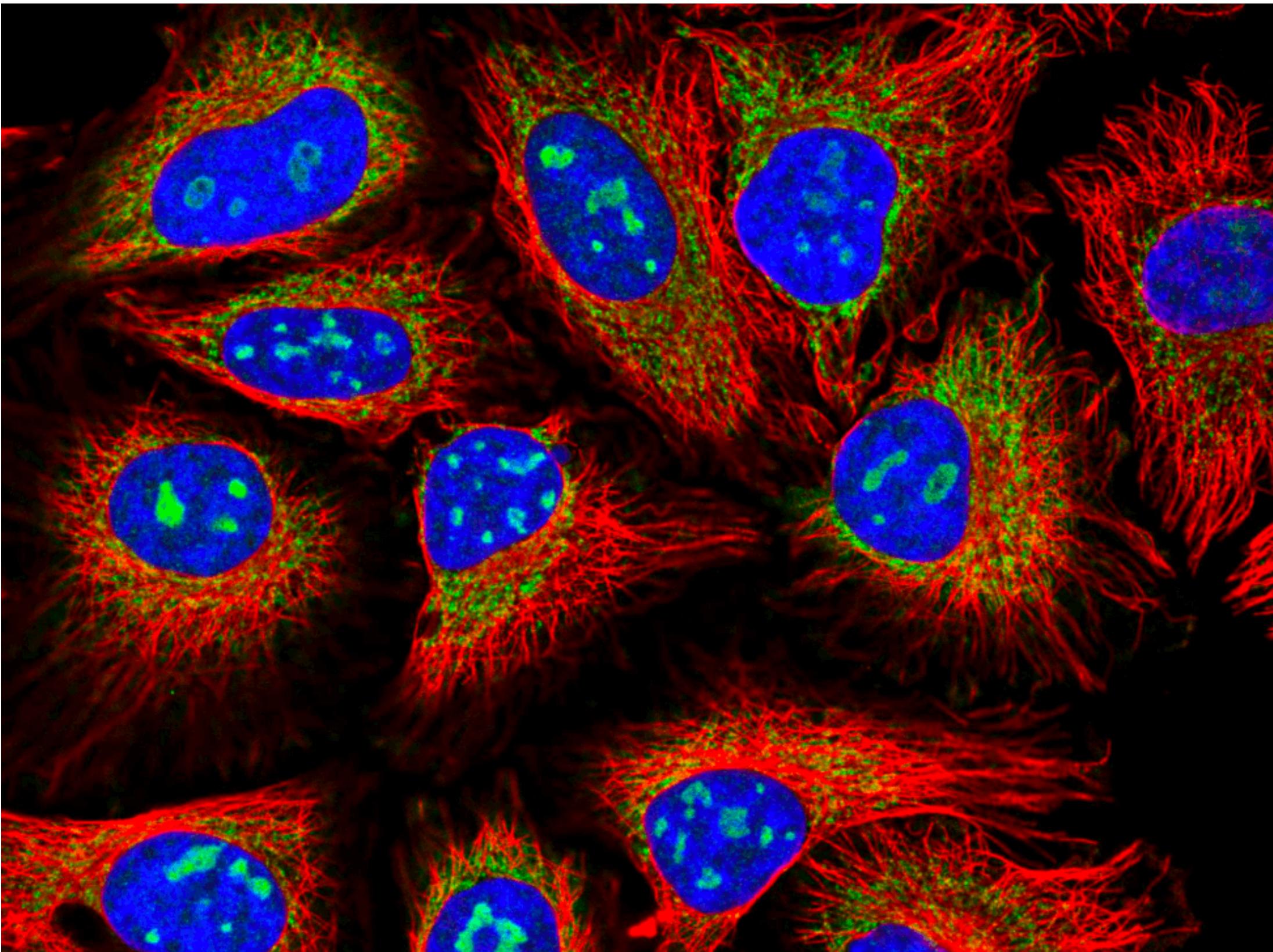
# Softmax

Turns outputs  $f$  into probabilities (sum up to 1 across  $k$  classes)



# Classification Tasks at Kaggle

Classify human protein microscope images into 28 categories

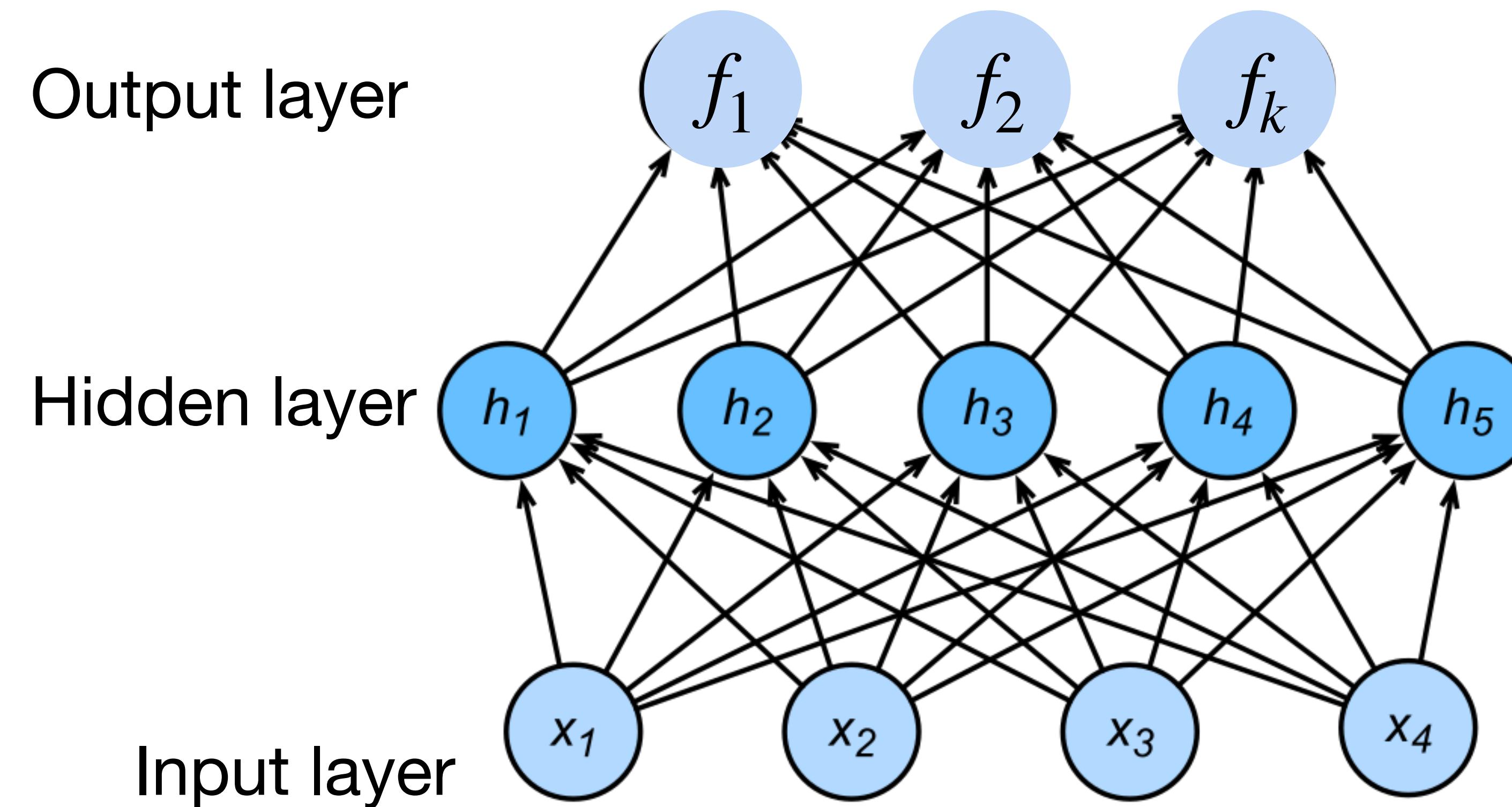


- 0. Nucleoplasm
- 1. Nuclear membrane
- 2. Nucleoli
- 3. Nucleoli fibrillar
- 4. Nuclear speckles
- 5. Nuclear bodies
- 6. Endoplasmic reticu
- 7. Golgi apparatus
- 8. Peroxisomes
- 9. Endosomes
- 10. Lysosomes
- 11. Intermediate fila
- 12. Actin filaments
- 13. Focal adhesion si
- 14. Microtubules
- 15. Microtubule ends
- 16. Cytokinetic brida

<https://www.kaggle.com/c/human-protein-atlas-image-classification>

# More complicated neural networks

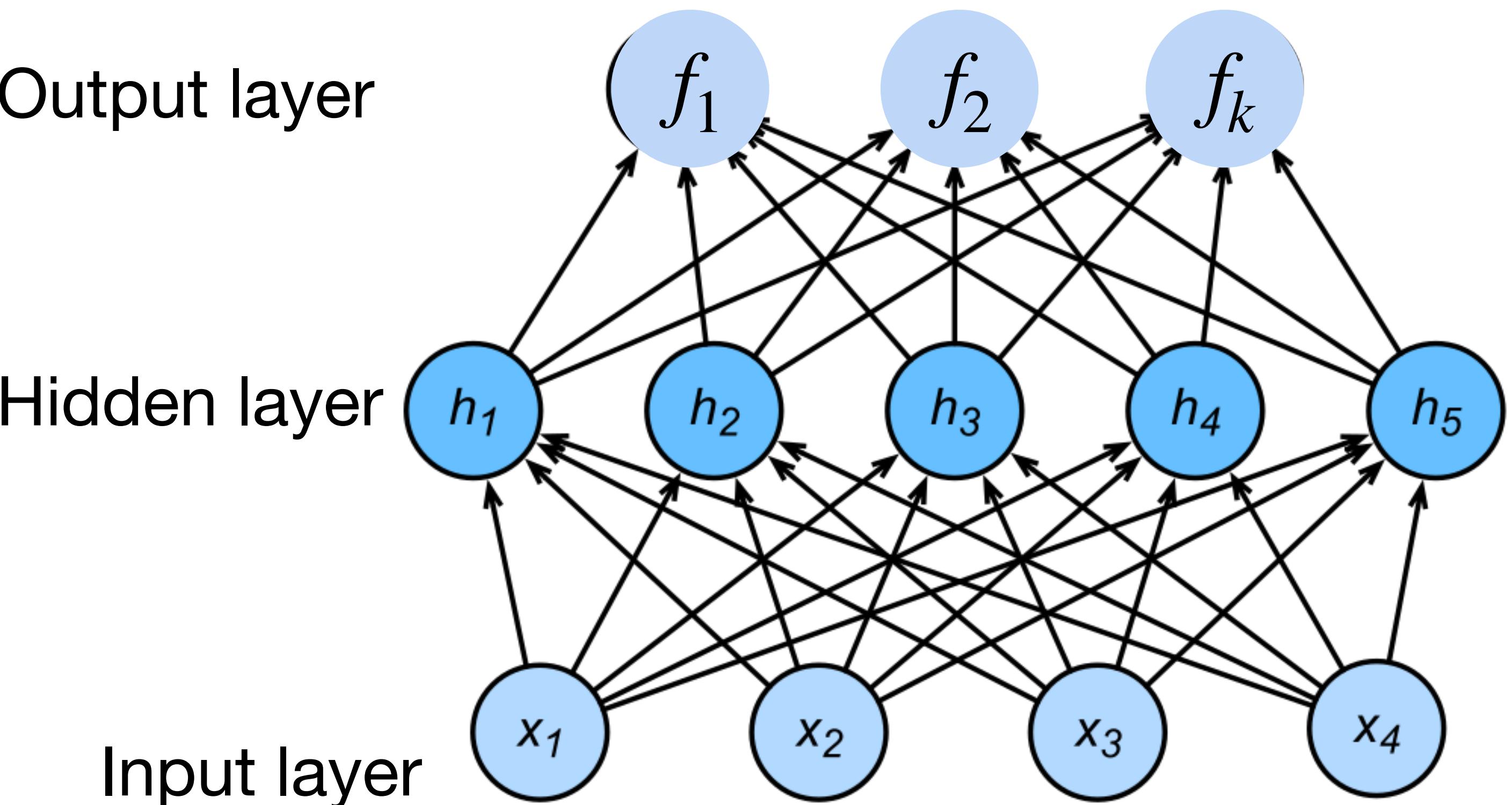
$$y_1, y_2, \dots, y_k = \text{softmax}(f_1, f_2, \dots, f_k)$$



# More complicated neural networks

- Input  $\mathbf{x} \in \mathbb{R}^d$
- Hidden layer  
 $\mathbf{W}^{(1)} \in \mathbb{R}^{m \times d}, \mathbf{b} \in \mathbb{R}^m$
- Outputs:  
 $\mathbf{h} = \sigma(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})$   
 $\mathbf{f} = \sigma(\mathbf{W}^{(2)}\mathbf{h} + \mathbf{b}^{(2)})$   
 $\mathbf{y} = \text{softmax}(\mathbf{f})$

$$y_1, y_2, \dots, y_k = \text{softmax}(f_1, f_2, \dots, f_k)$$



# More complicated neural networks: 3 hidden layers

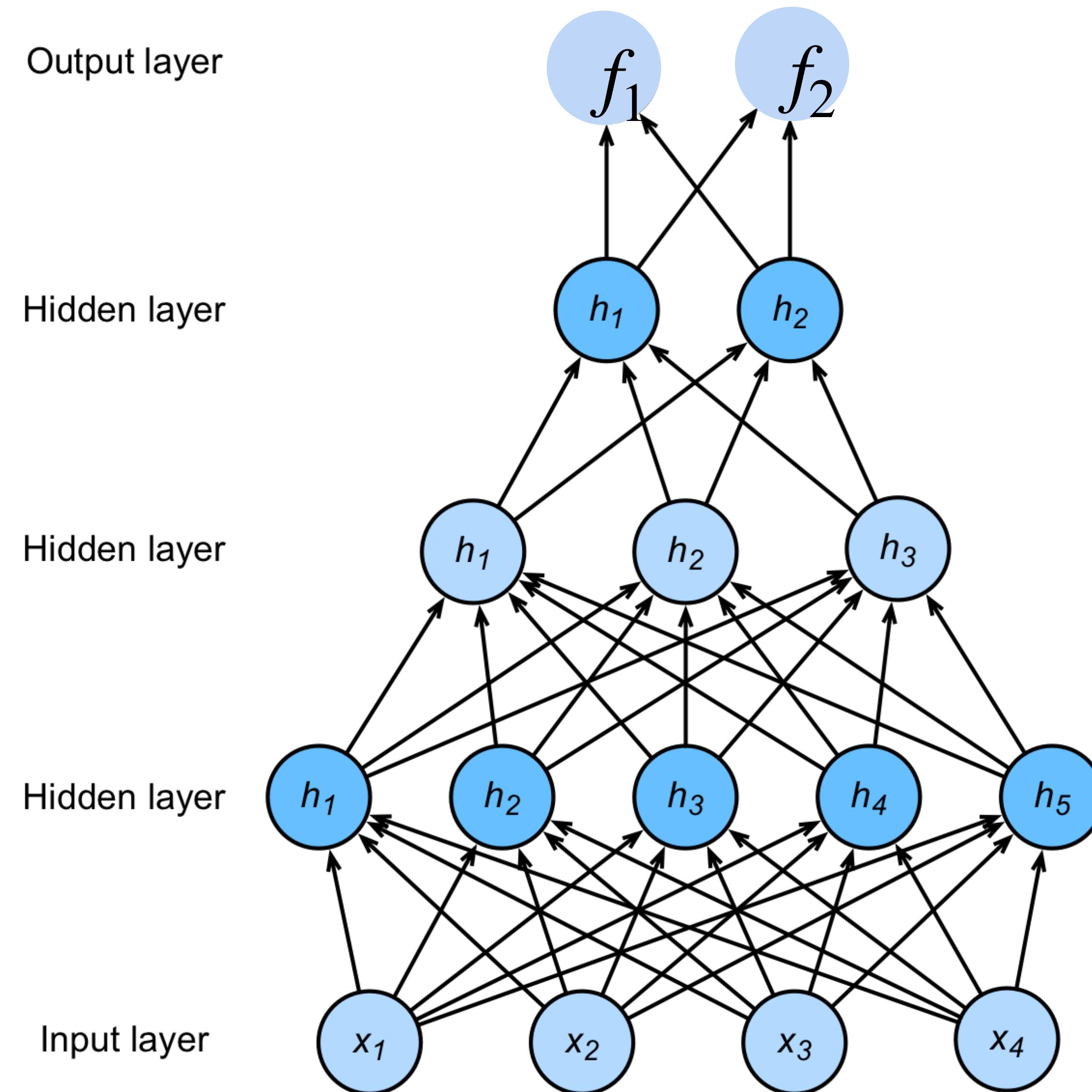
$$\mathbf{h}_1 = \sigma(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})$$

$$\mathbf{h}_2 = \sigma(\mathbf{W}^{(2)}\mathbf{h}_1 + \mathbf{b}^{(2)})$$

$$\mathbf{h}_3 = \sigma(\mathbf{W}^{(3)}\mathbf{h}_2 + \mathbf{b}^{(3)})$$

$$\mathbf{f} = \mathbf{W}^{(4)}\mathbf{h}_3 + \mathbf{b}^{(4)}$$

$$\mathbf{y} = \text{softmax}(\mathbf{f})$$



# Quiz Break

Which output function is often used for multi-class classification tasks?

- A Sigmoid function
- B Rectified Linear Unit (ReLU)
- C Softmax function
- D Max function

# Quiz Break

Suppose you are given a 3-layer multilayer perceptron (2 hidden layers  $h_1$  and  $h_2$  and 1 output layer). All activation functions are sigmoids, and the output layer uses a softmax function. Suppose  $h_1$  has 1024 units and  $h_2$  has 512 units. Given a dataset with 2 input features and 3 unique class labels, how many learnable parameters does the perceptron have in total?

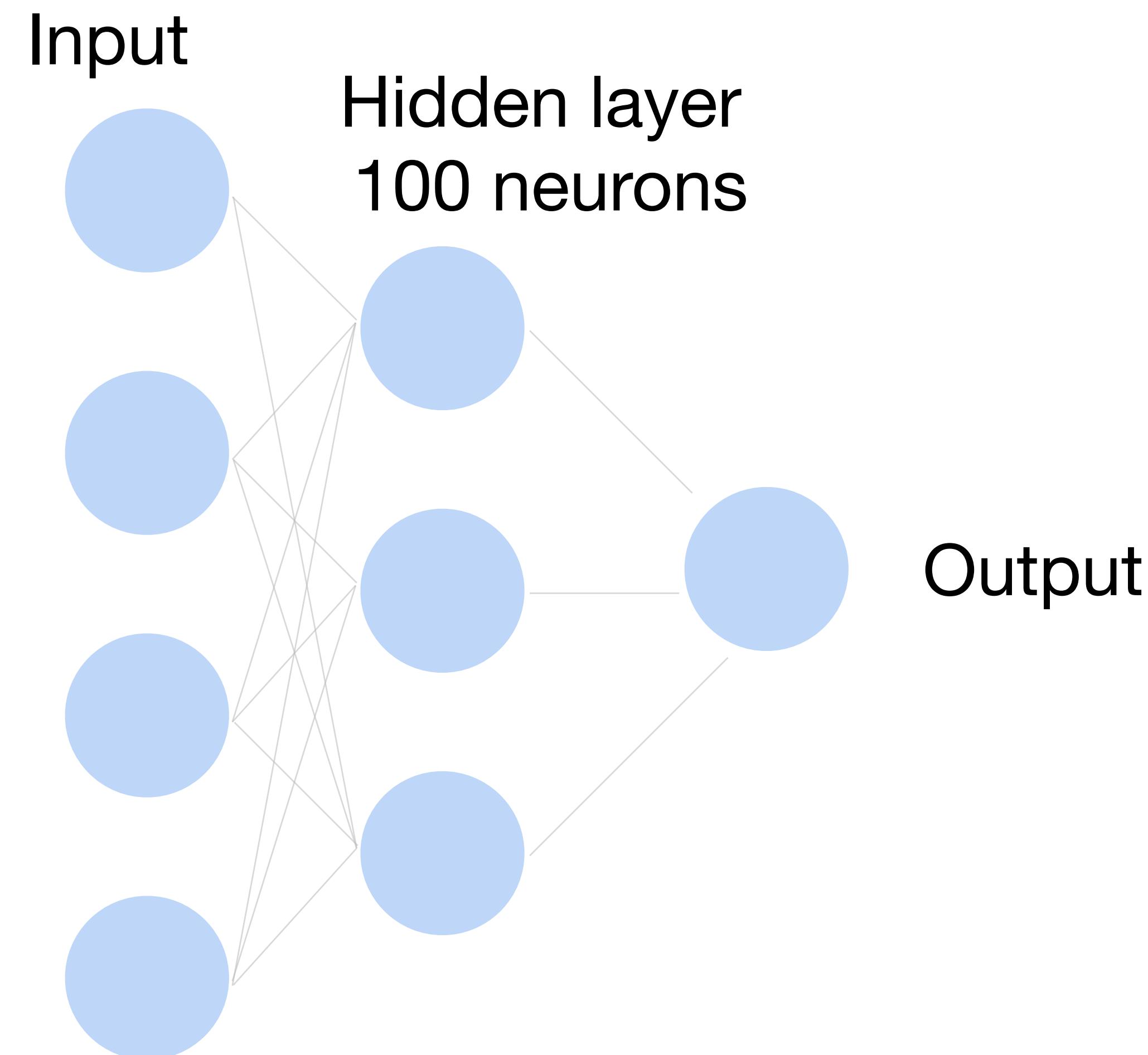
# Quiz Break

Consider a three-layer network with **linear Perceptrons** for binary classification. The hidden layer has 3 neurons. Can the network represent a XOR problem?

- a) Yes
- b) No

# How to train a neural network?

**Classify cats vs. dogs**



# How to train a neural network? 2-class

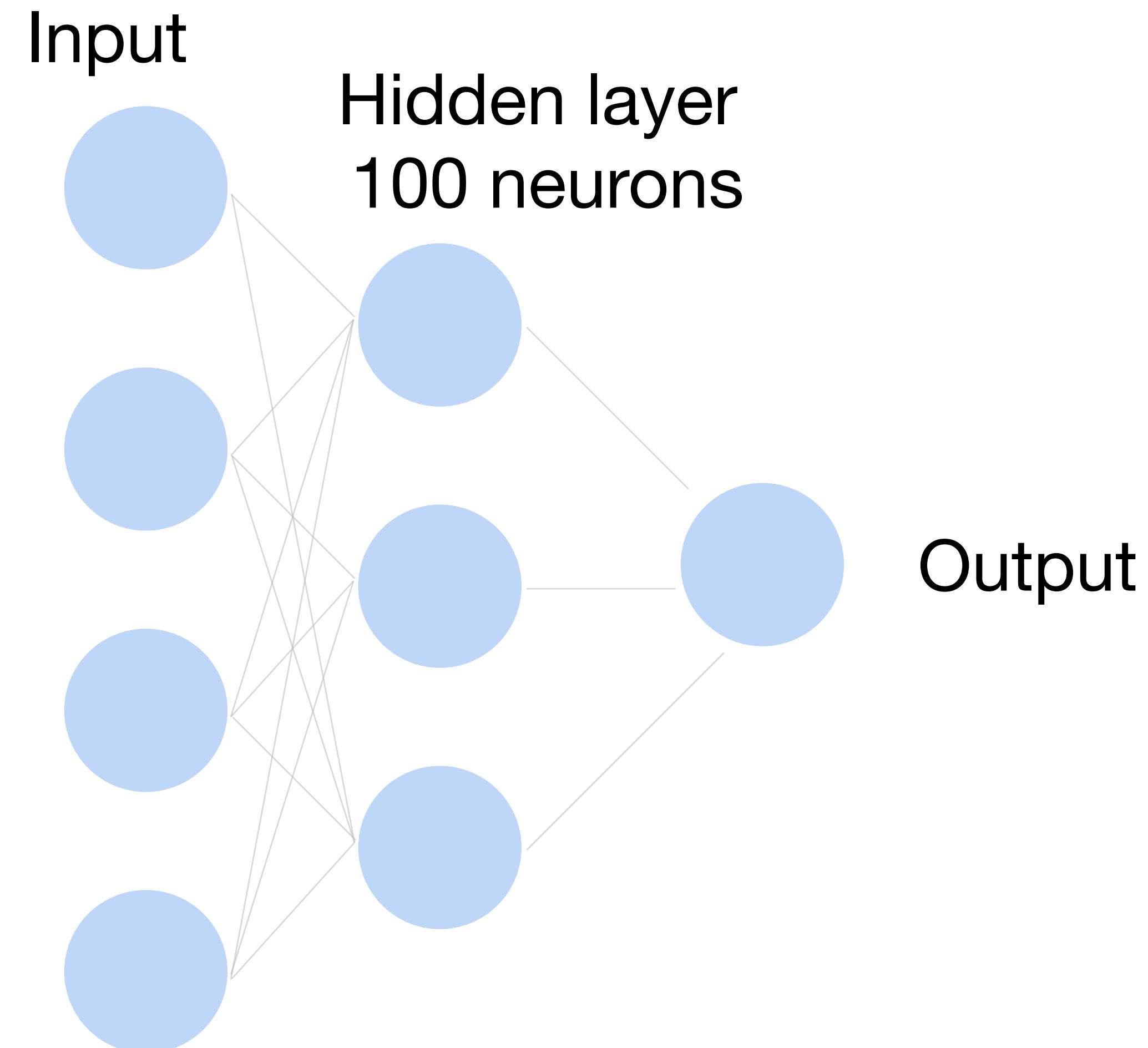
$\mathbf{x}_i \in \mathbb{R}^d$  : One training data point in the training set D

$\hat{y}_i$  Model output for  $\mathbf{x}_i$

$y_i$  Ground truth label for  $\mathbf{x}_i$

**Learning by matching the output to the label**

We want  $\hat{y}_i \rightarrow 1$  when  $y_i = 1$ ,  
and  $\hat{y}_i \rightarrow 0$  when  $y_i = 0$

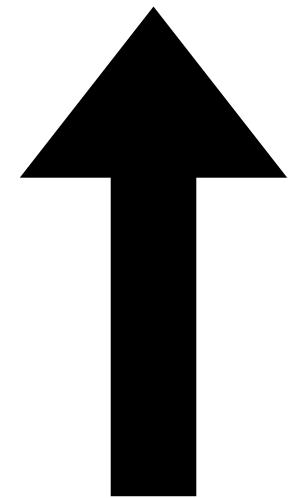


# How to train a neural network? 2-class

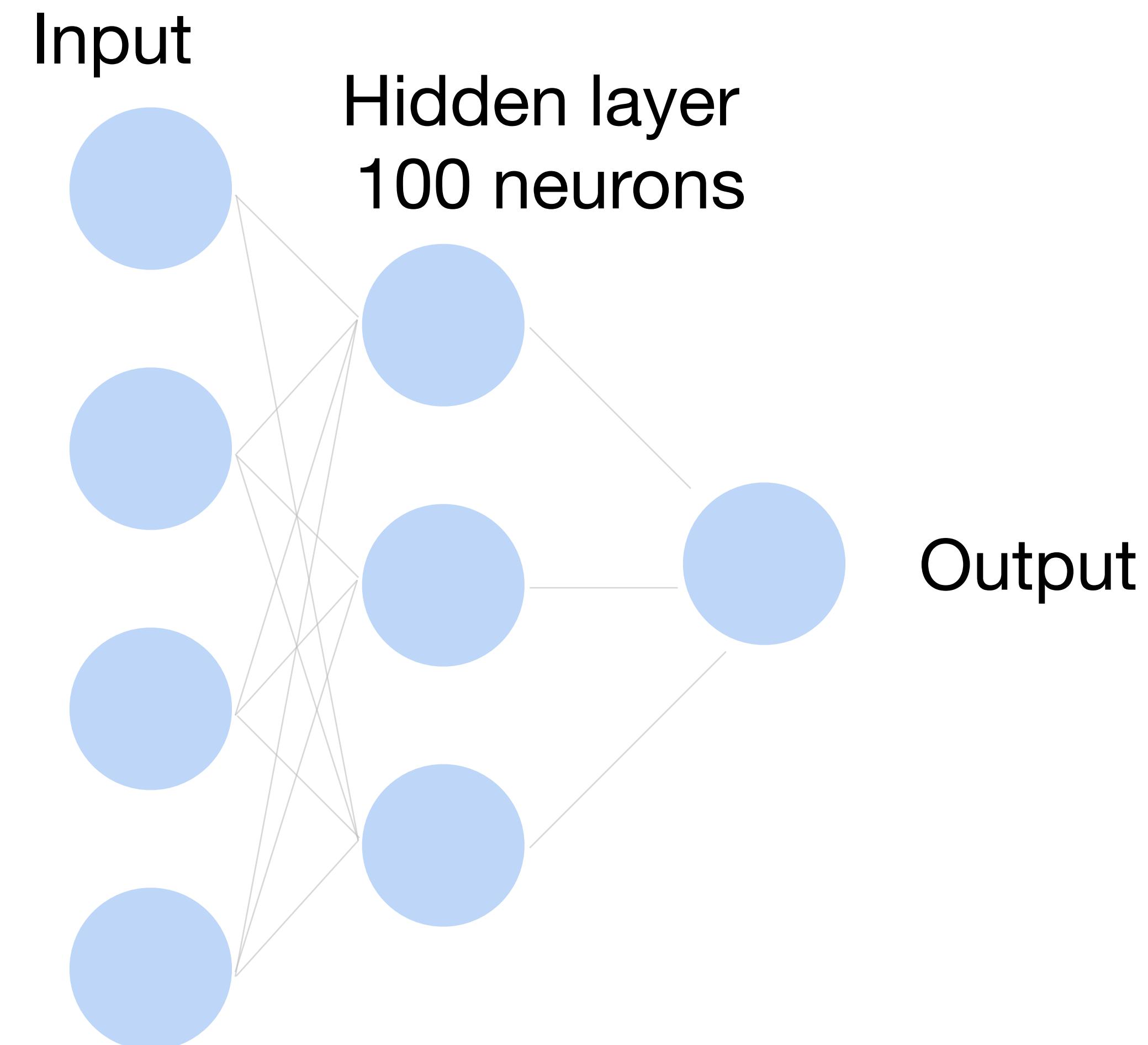
Loss function:  $\frac{1}{|D|} \sum_i \ell(\mathbf{x}_i, y_i)$

Per-sample loss:

$$\ell(\mathbf{x}_i, y_i) = -y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)$$



Also known as **binary cross-entropy loss**



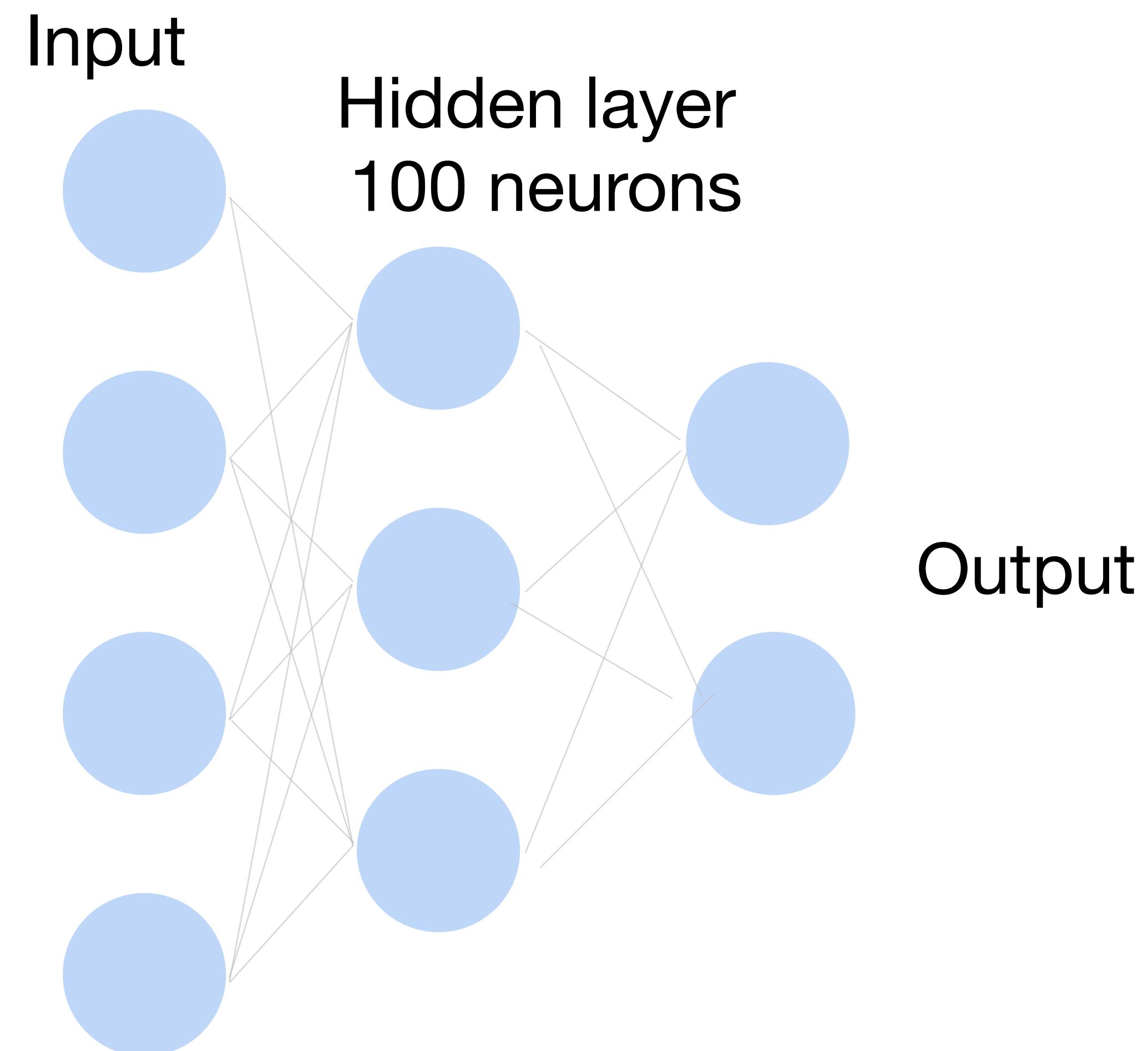
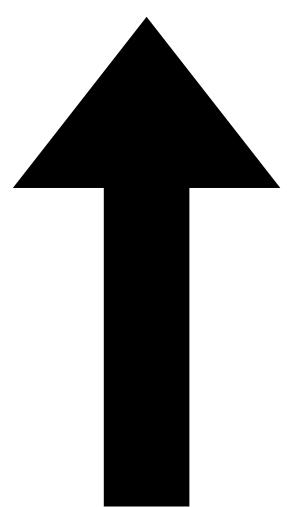
# How to train a neural network? $K$ -class

Loss function:  $\frac{1}{|D|} \sum_i \ell(\mathbf{x}_i, y_i)$

Per-sample loss:

$$\ell(\mathbf{x}_i, y_i) = \sum_{j=1}^K -y_{ij} \log \hat{p}_{ij}$$

Also known as **cross-entropy loss**  
or **softmax loss**

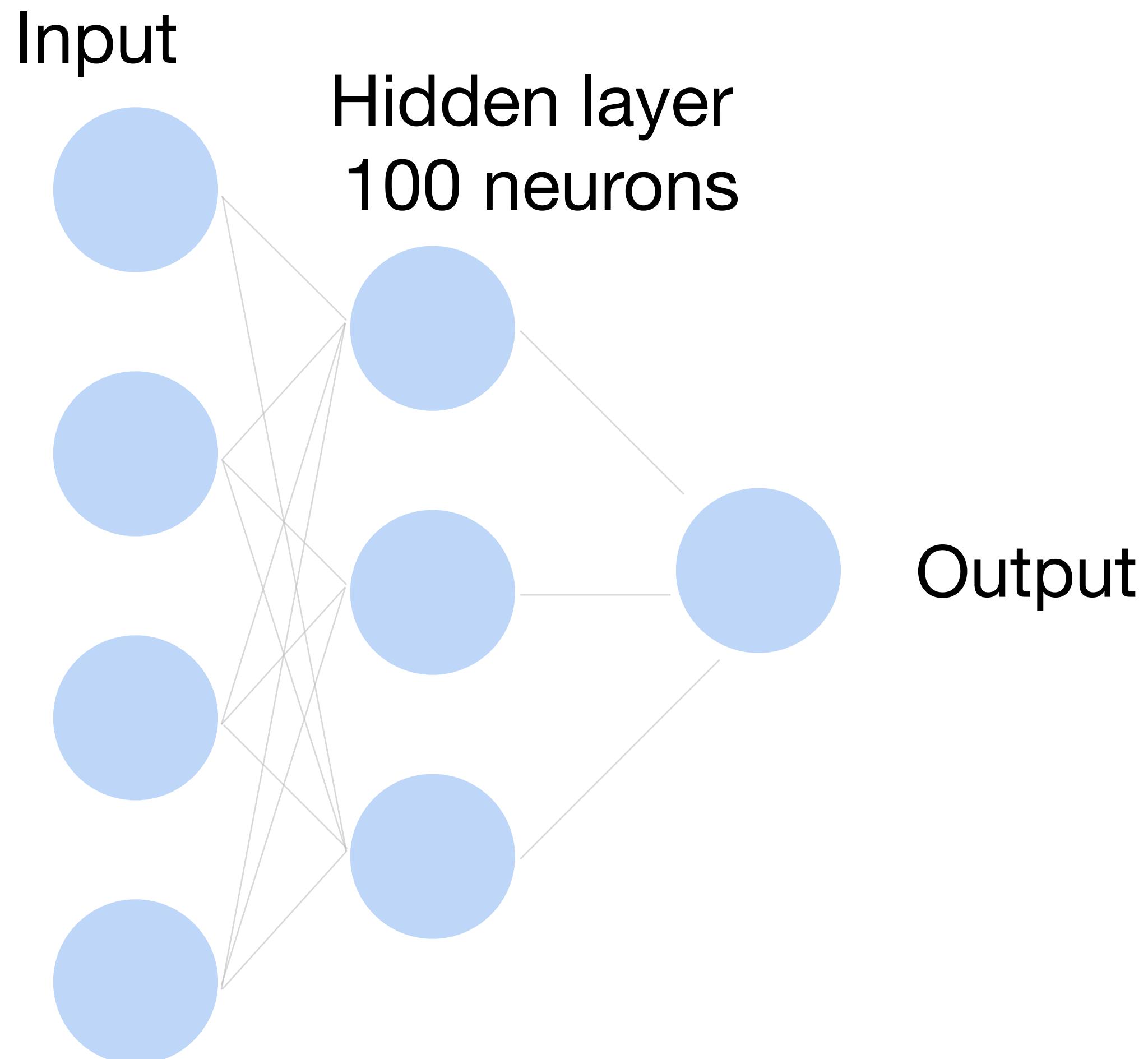


# How to train a neural network?

Update the weights  $W$  to minimize the loss function

$$L = \frac{1}{|D|} \sum_i \ell(\mathbf{x}_i, y_i)$$

Use gradient descent!



# Gradient Descent

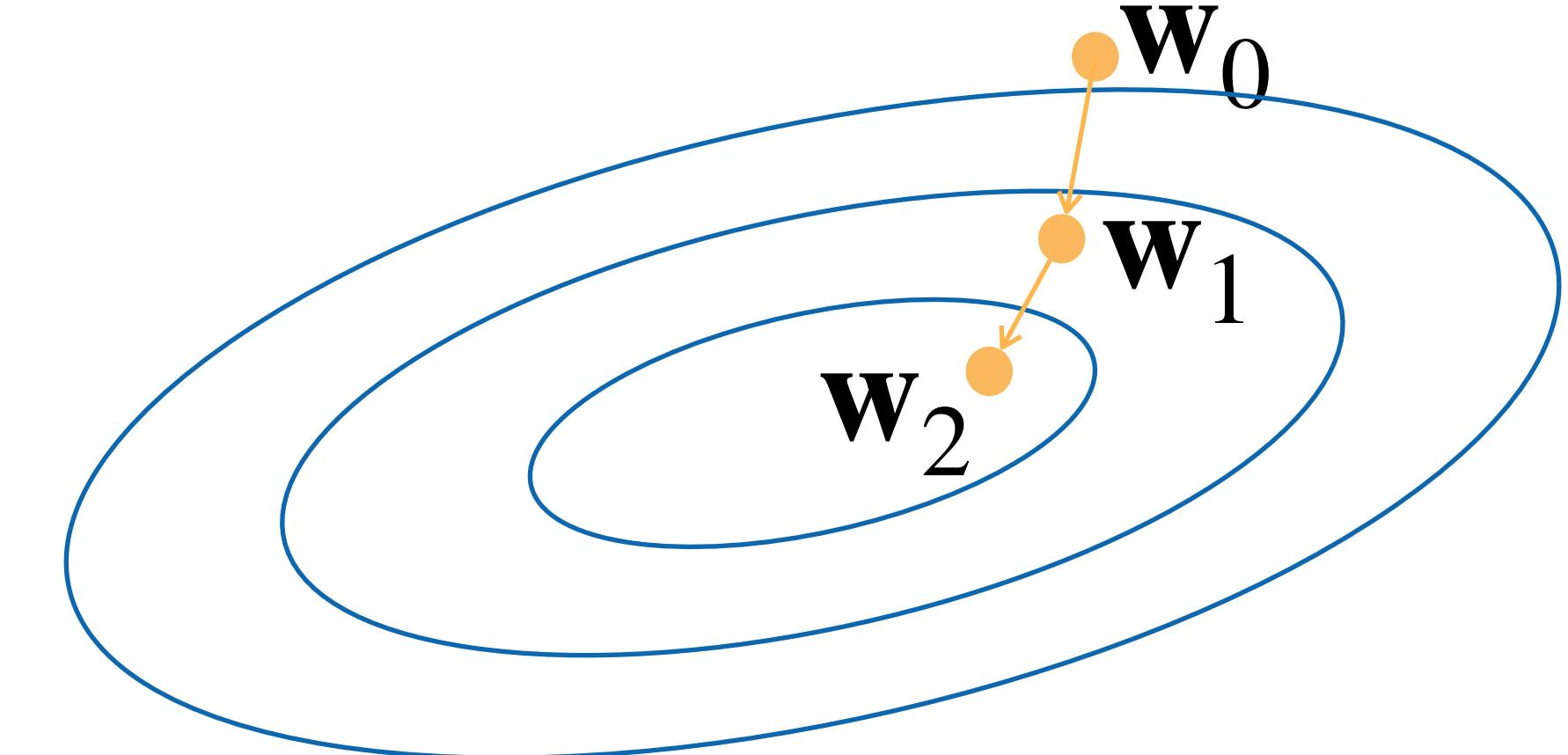
- Choose a learning rate  $\alpha > 0$
- Initialize the model parameters  $\mathbf{w}_0$
- For  $t = 1, 2, \dots$

- Update parameters:

$$\begin{aligned}\mathbf{w}_t &= \mathbf{w}_{t-1} - \alpha \frac{\partial L}{\partial \mathbf{w}_{t-1}} \\ &= \mathbf{w}_{t-1} - \alpha \frac{1}{|D|} \sum_{i \in D} \frac{\partial \ell(\mathbf{x}_i, y_i)}{\partial \mathbf{w}_{t-1}}\end{aligned}$$

D can  
be very large.  
Expensive

- Repeat until converges



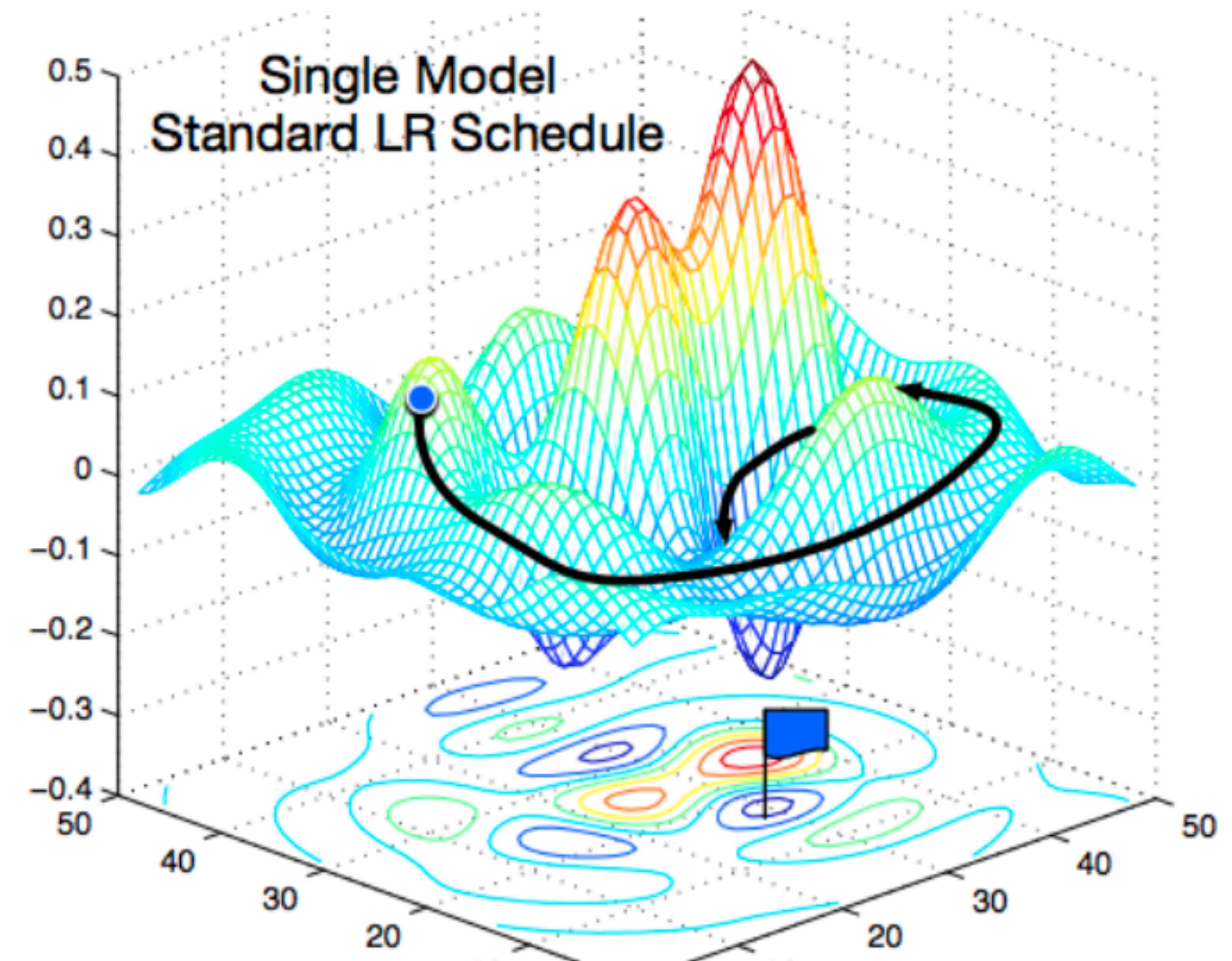
# Minibatch Stochastic Gradient Descent

- Choose a learning rate  $\alpha > 0$
- Initialize the model parameters  $\mathbf{w}_0$
- For  $t = 1, 2, \dots$ 
  - **Randomly sample a subset (mini-batch)  $\hat{D} \subset D$**   
Update parameters:

$$\mathbf{w}_t = \mathbf{w}_{t-1} - \alpha \frac{1}{|\hat{D}|} \sum_{i \in \hat{D}} \frac{\partial \ell(\mathbf{x}_i, y_i)}{\partial \mathbf{w}_{t-1}}$$

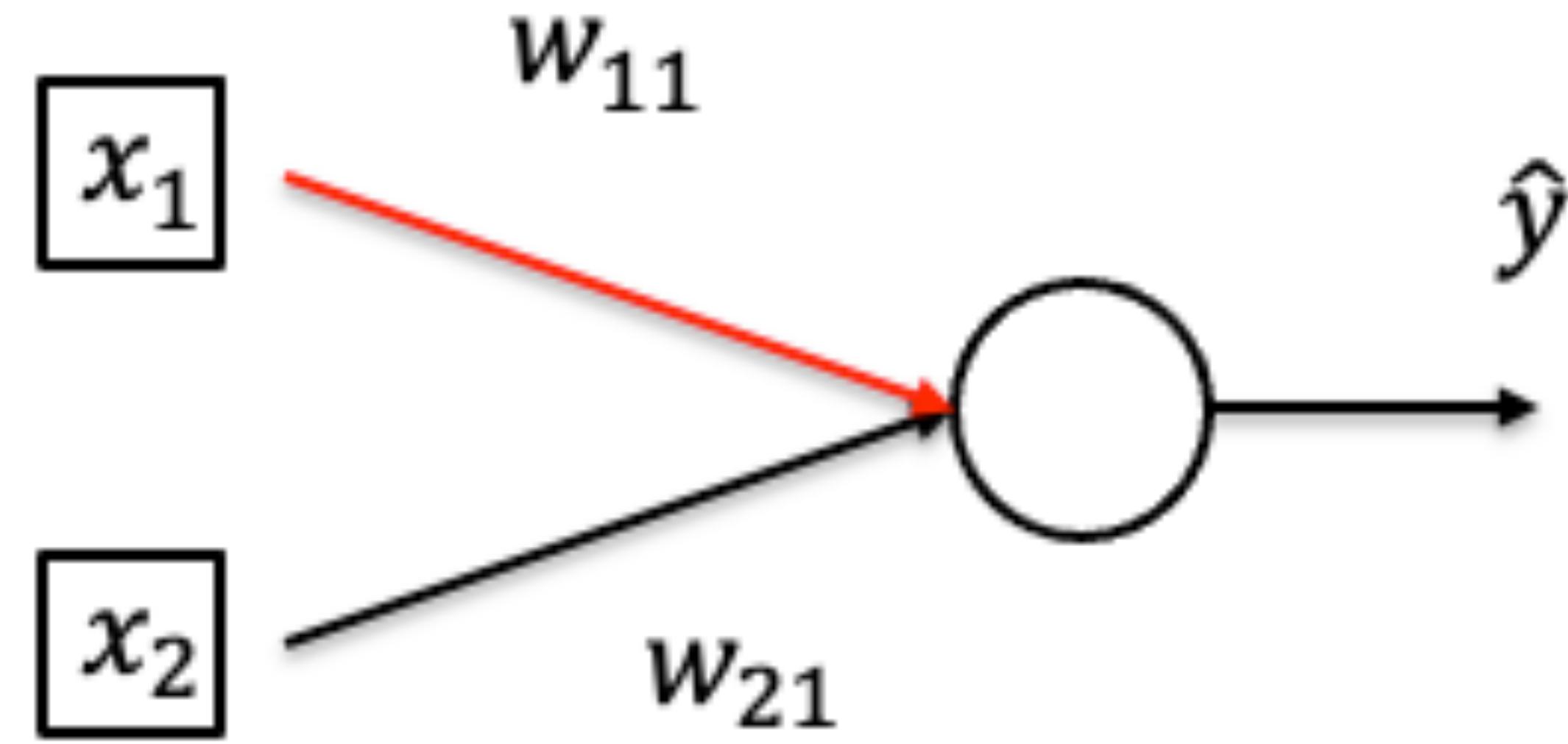
- Repeat until converges

# Non-convex Optimization



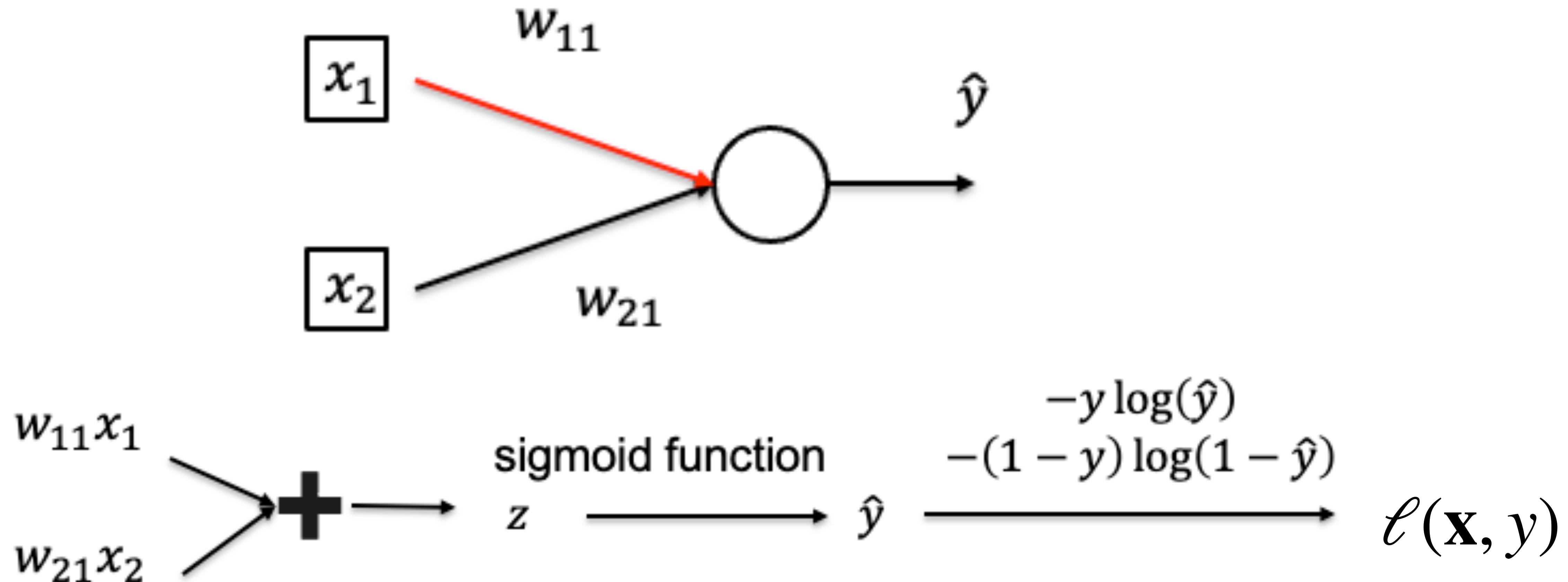
[Gao and Li et al., 2018]

# Calculate Gradient (on one data point)

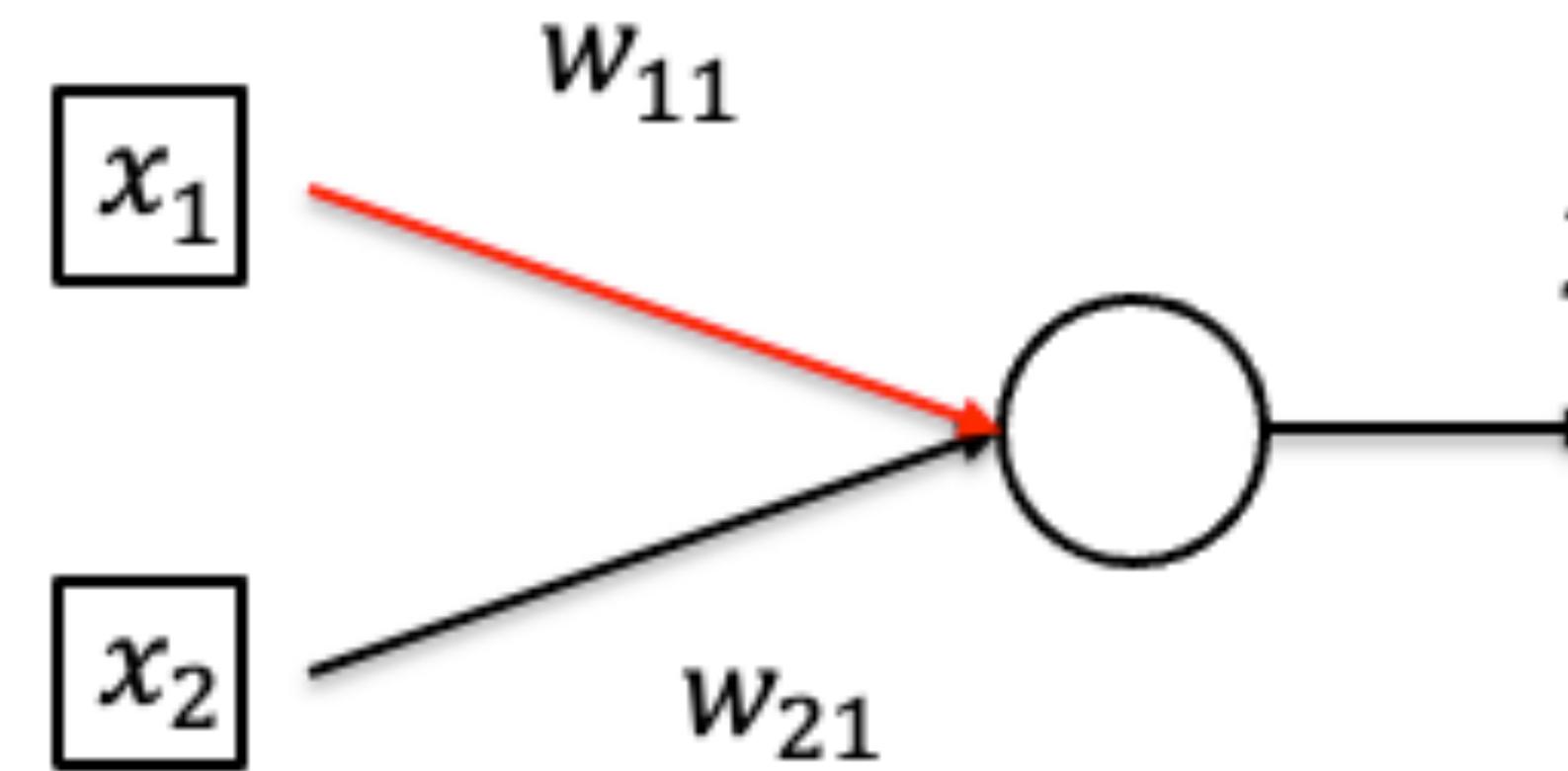


- Want to compute  $\frac{\partial \ell(\mathbf{x}, y)}{\partial w_{11}}$

# Calculate Gradient (on one data point)



# Calculate Gradient (on one data point)



A computational graph illustrating the forward pass and backpropagation for a logistic regression model. The forward pass consists of three main steps:

- Inputs  $w_{11}x_1$  and  $w_{21}x_2$  are summed to produce  $z$ .
- $z$  is passed through a sigmoid function to produce  $\hat{y}$ .
- $\hat{y}$  is used to calculate the loss function  $\ell(\mathbf{x}, y)$ .

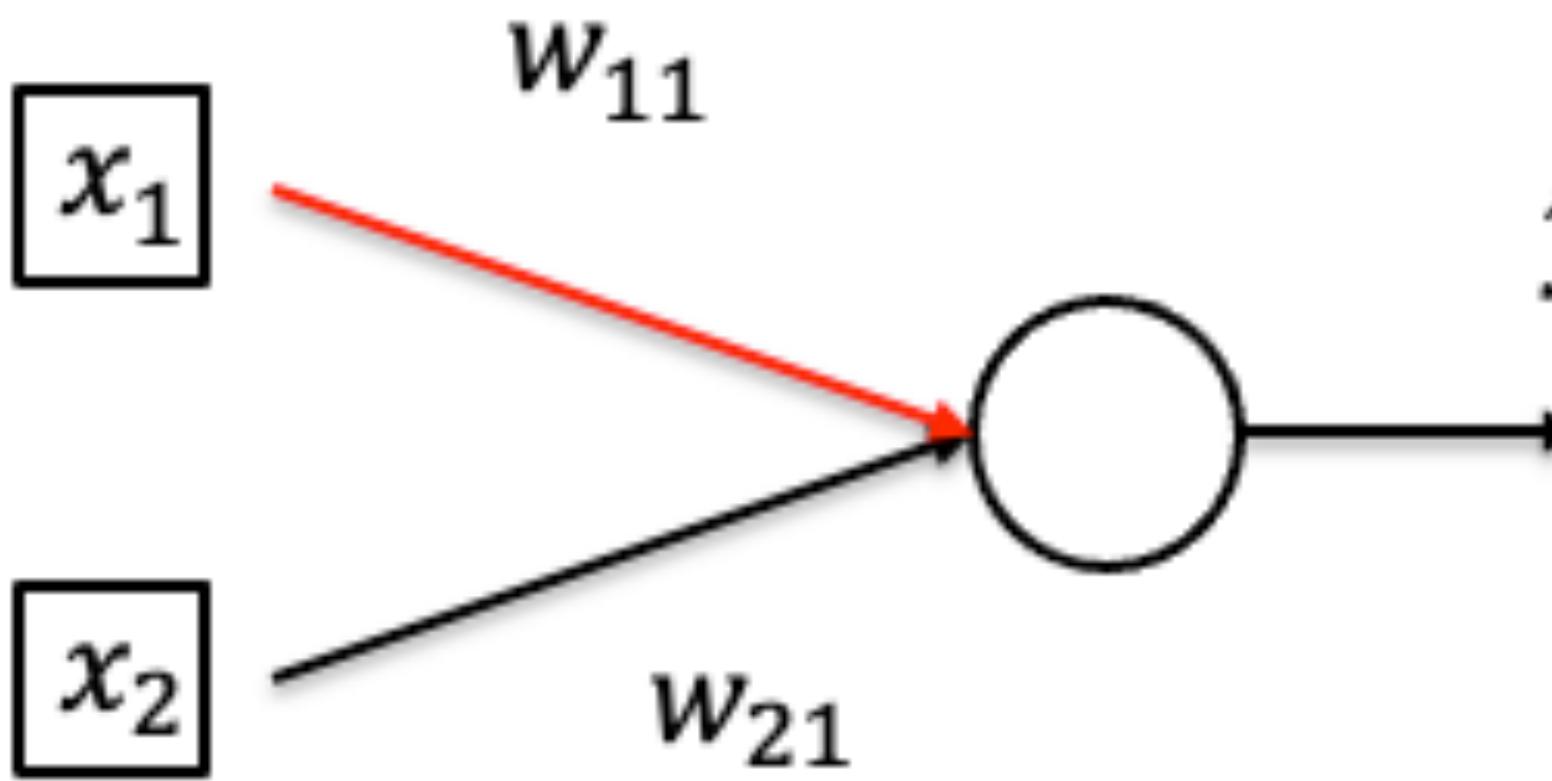
The backpropagation step shows the derivative of the loss function with respect to  $\hat{y}$ :

$$\frac{\partial \ell(\mathbf{x}, y)}{\partial \hat{y}} = \frac{-y \log(\hat{y})}{-(1 - y) \log(1 - \hat{y})} = \frac{1 - y}{1 - \hat{y}} - \frac{y}{\hat{y}}$$

- By chain rule:

$$\frac{\partial l}{\partial w_{11}} = \frac{\partial l}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial w_{11}}$$

# Calculate Gradient (on one data point)

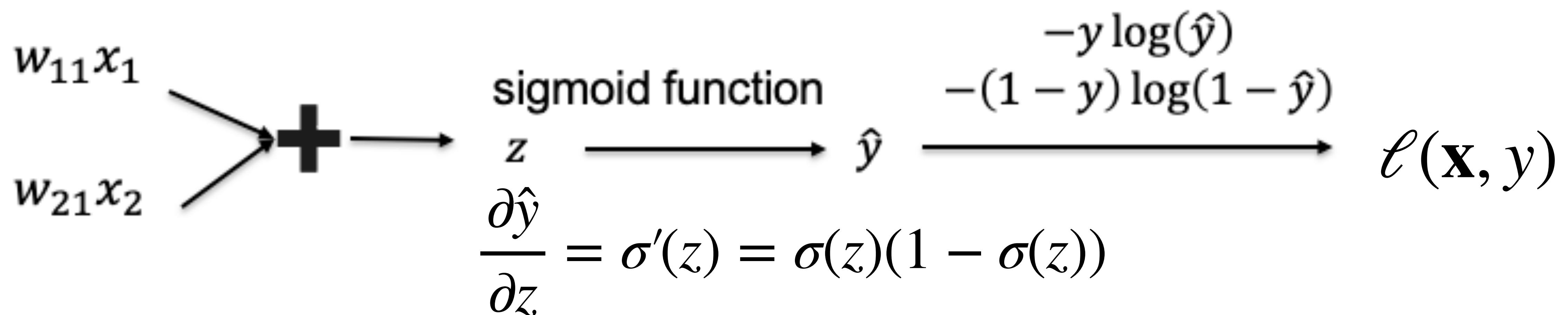
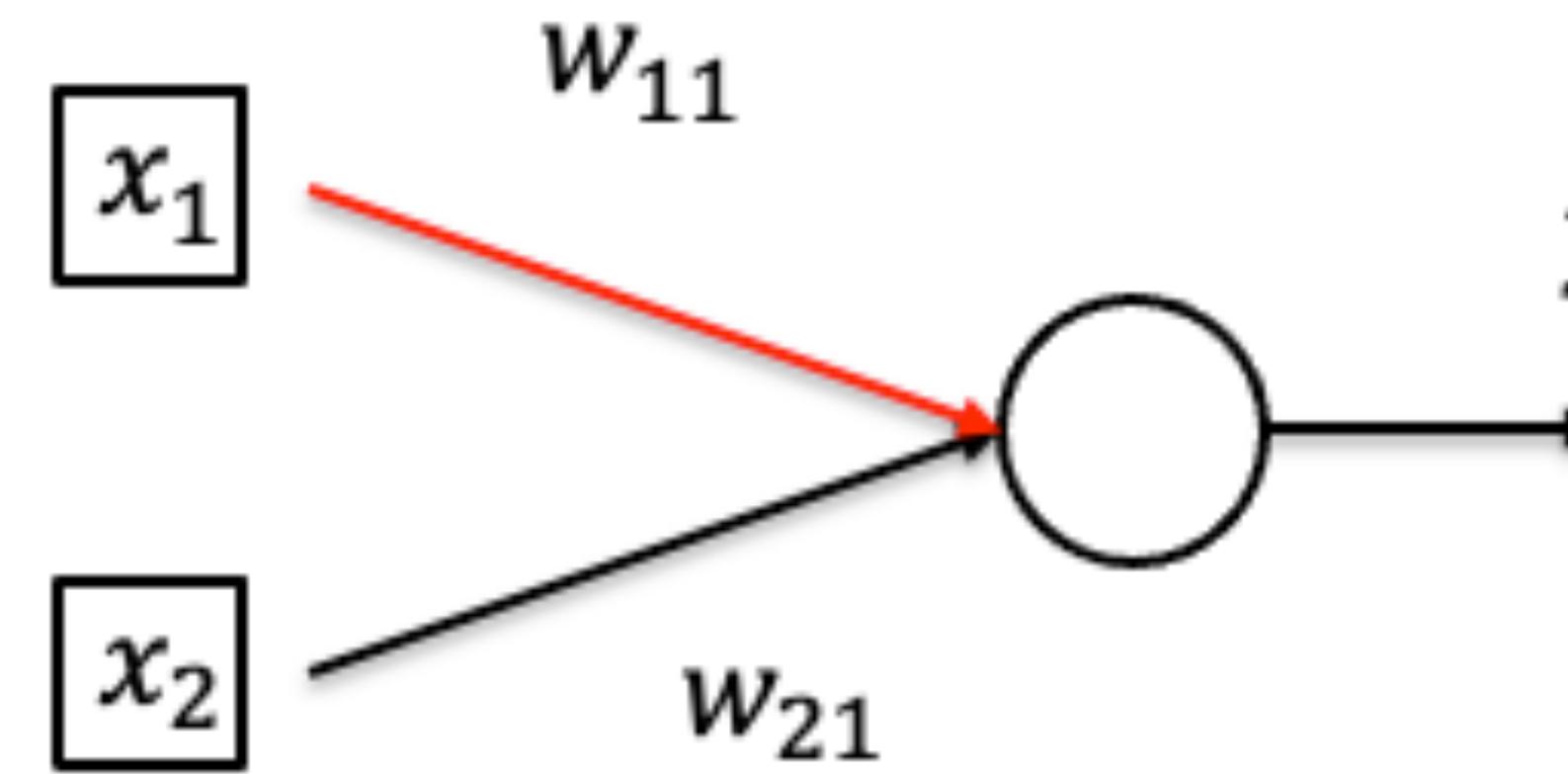


$$\begin{array}{ccccc} w_{11}x_1 & \xrightarrow{\text{+}} & z & \xrightarrow{\text{sigmoid function}} & \hat{y} \\ w_{21}x_2 & & & & \xrightarrow{-y \log(\hat{y})} \\ & & & & -(1 - \hat{y}) \log(1 - \hat{y}) \\ & & & & \ell(\mathbf{x}, y) \\ \frac{\partial \hat{y}}{\partial z} = \sigma'(z) & & & \frac{\partial \ell(\mathbf{x}, y)}{\partial \hat{y}} = \frac{1 - y}{1 - \hat{y}} - \frac{y}{\hat{y}} & \end{array}$$

- By chain rule:

$$\frac{\partial l}{\partial w_{11}} = \frac{\partial l}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} x_1$$

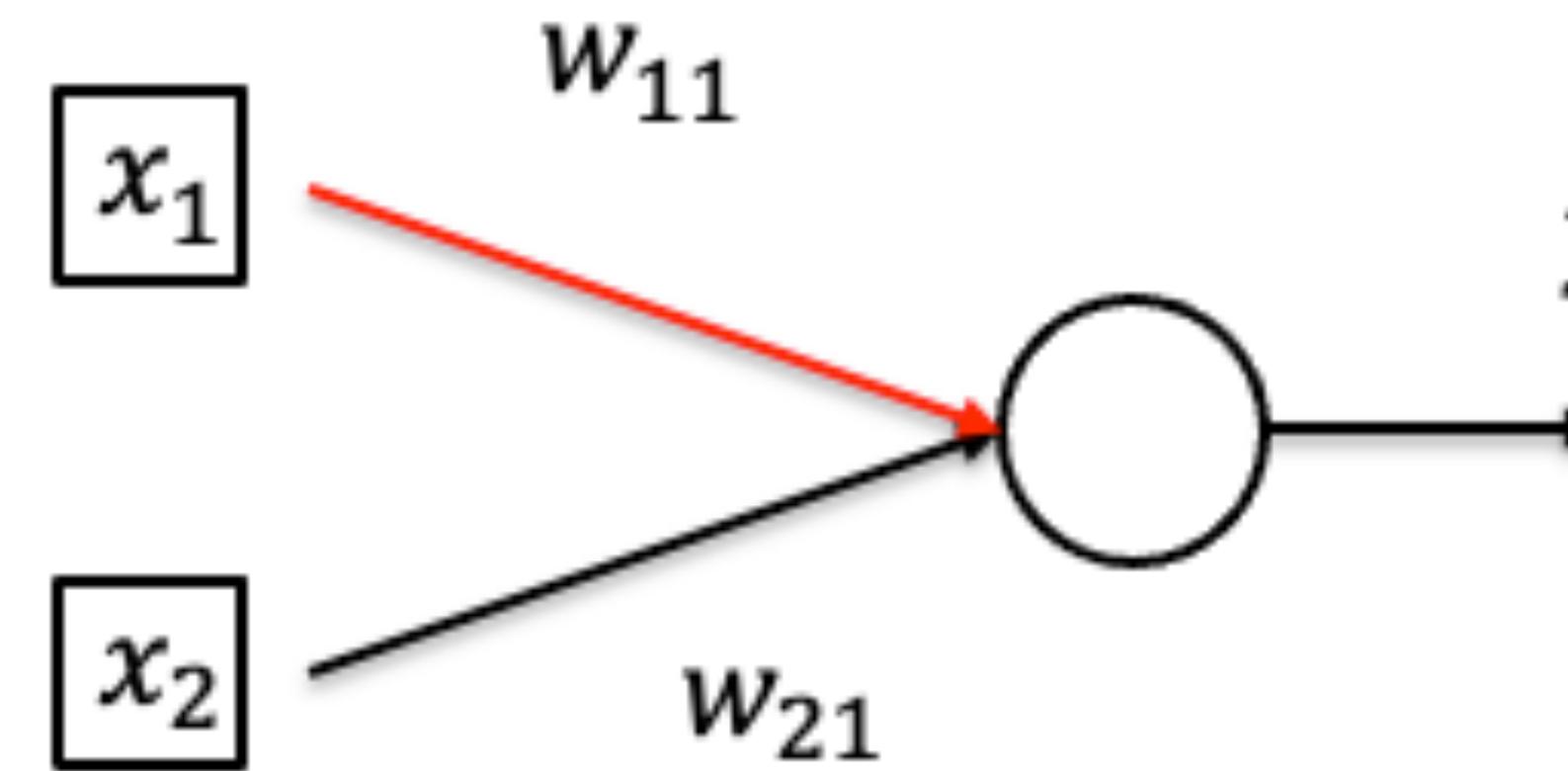
# Calculate Gradient (on one data point)



- By chain rule:

$$\frac{\partial l}{\partial w_{11}} = \frac{\partial l}{\partial \hat{y}} \hat{y}(1 - \hat{y})x_1$$

# Calculate Gradient (on one data point)



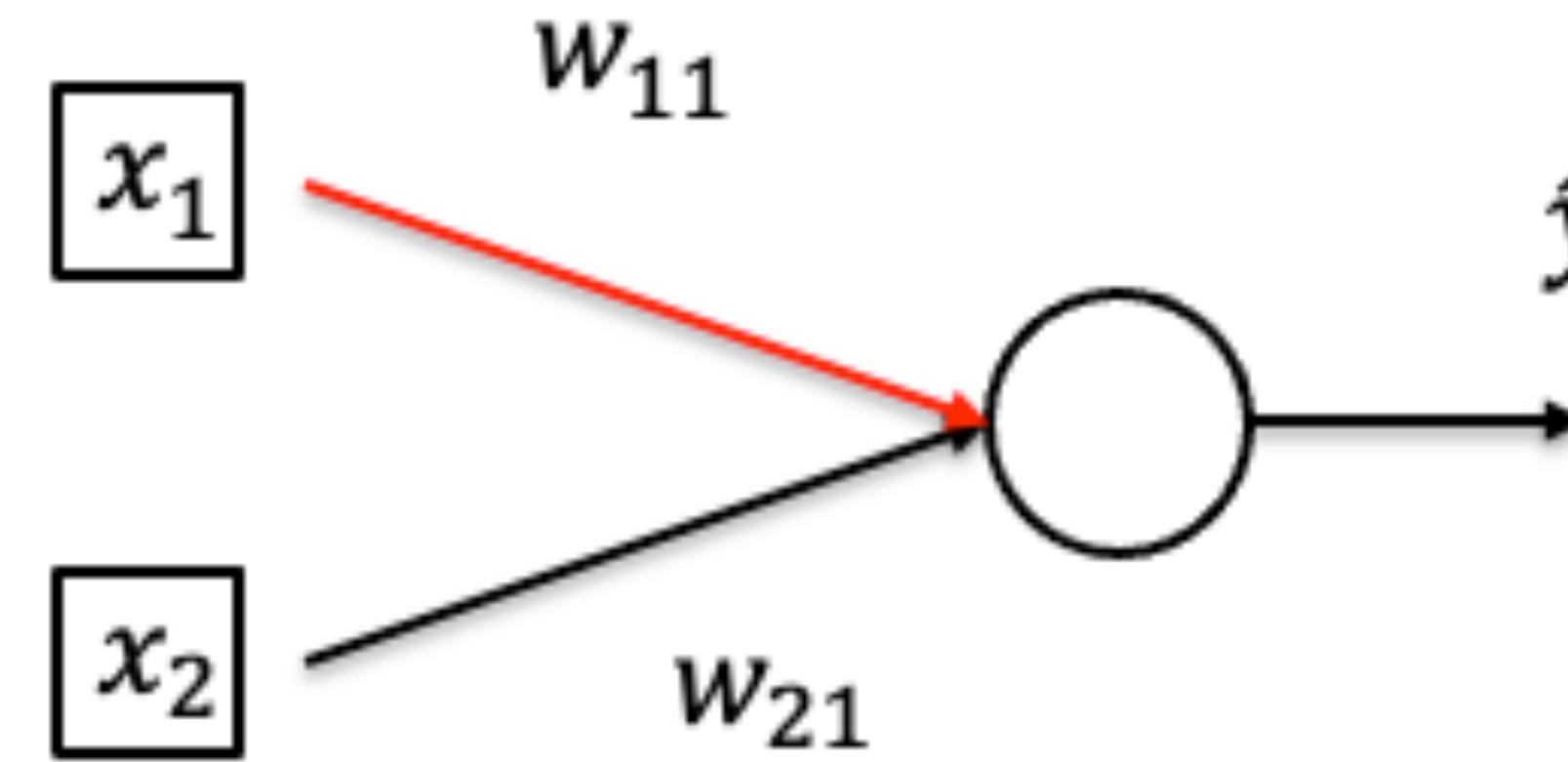
The diagram illustrates the calculation of the loss function  $\ell(\mathbf{x}, y)$  from the inputs  $x_1$  and  $x_2$ . The inputs are multiplied by their respective weights  $w_{11}$  and  $w_{21}$ , then summed. This sum is passed through a sigmoid function to produce the predicted output  $\hat{y}$ . The loss function is then calculated as  $-(y \log(\hat{y}) + (1 - y) \log(1 - \hat{y}))$ .

$$\begin{array}{c} w_{11}x_1 \\ w_{21}x_2 \\ \hline + \end{array} \quad \text{sigmoid function} \quad z \longrightarrow \hat{y} \quad \frac{-y \log(\hat{y})}{-(1 - y) \log(1 - \hat{y})} \longrightarrow \ell(\mathbf{x}, y)$$
$$\frac{\partial \hat{y}}{\partial z} = \sigma'(z) = \sigma(z)(1 - \sigma(z))$$

- By chain rule:

$$\frac{\partial l}{\partial w_{11}} = \left( \frac{1 - y}{1 - \hat{y}} - \frac{y}{\hat{y}} \right) \hat{y}(1 - \hat{y})x_1$$

# Calculate Gradient (on one data point)



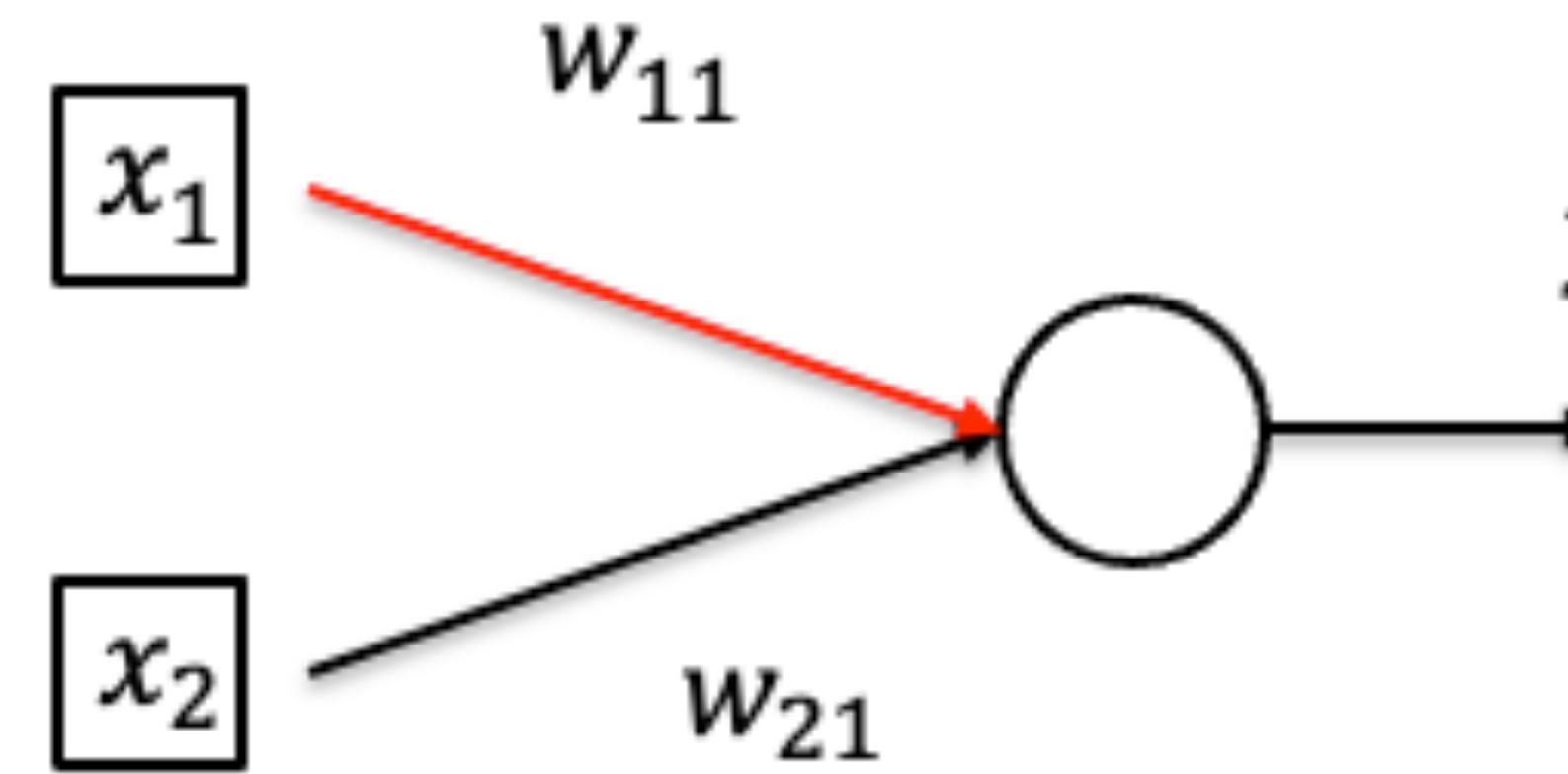
A computational graph illustrating the forward pass of a neural network. It starts with two inputs,  $w_{11}x_1$  and  $w_{21}x_2$ , which are summed at a plus sign node. The result is passed through a "sigmoid function" node, which outputs  $\hat{y}$ . Finally, the output  $\hat{y}$  is used to calculate the cross-entropy loss,  $\ell(\mathbf{x}, y)$ . Below the graph, the derivative of the sigmoid function is given as  $\frac{\partial \hat{y}}{\partial z} = \sigma'(z) = \sigma(z)(1 - \sigma(z))$ .

$$\begin{array}{ccccc} w_{11}x_1 & + & \text{sigmoid function} & -y \log(\hat{y}) & \ell(\mathbf{x}, y) \\ w_{21}x_2 & & z \longrightarrow \hat{y} & -(1 - \hat{y}) \log(1 - \hat{y}) & \\ & & & \frac{\partial \hat{y}}{\partial z} = \sigma'(z) = \sigma(z)(1 - \sigma(z)) & \end{array}$$

- By chain rule:

$$\frac{\partial l}{\partial w_{11}} = (\hat{y} - y)x_1$$

# Calculate Gradient (on one data point)



A computational graph illustrating the forward pass and gradients for a single data point. The forward pass consists of three main steps:

- Input features  $x_1$  and  $x_2$  are multiplied by weights  $w_{11}$  and  $w_{21}$  respectively, and then summed.
- The result of the summation is passed through a sigmoid function to produce the predicted output  $\hat{y}$ .
- The loss function  $\ell(\mathbf{x}, y)$  is calculated based on the predicted output  $\hat{y}$  and the true target  $y$ .

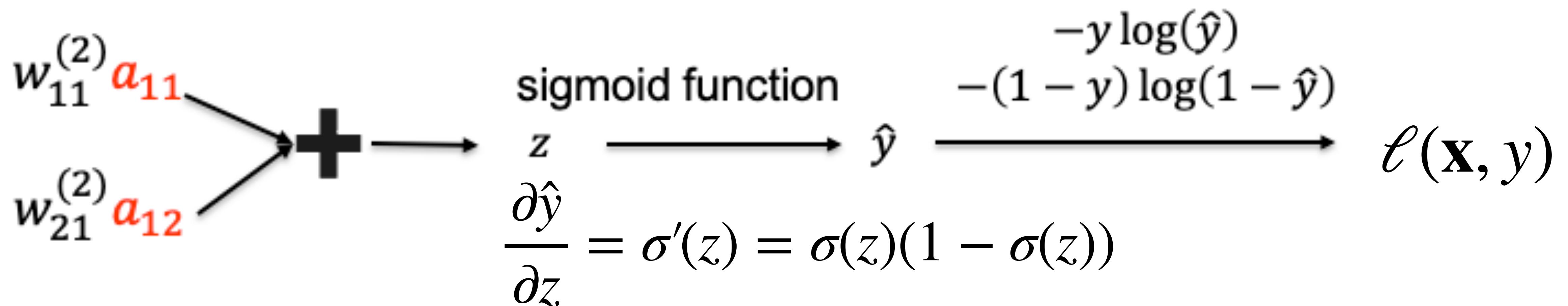
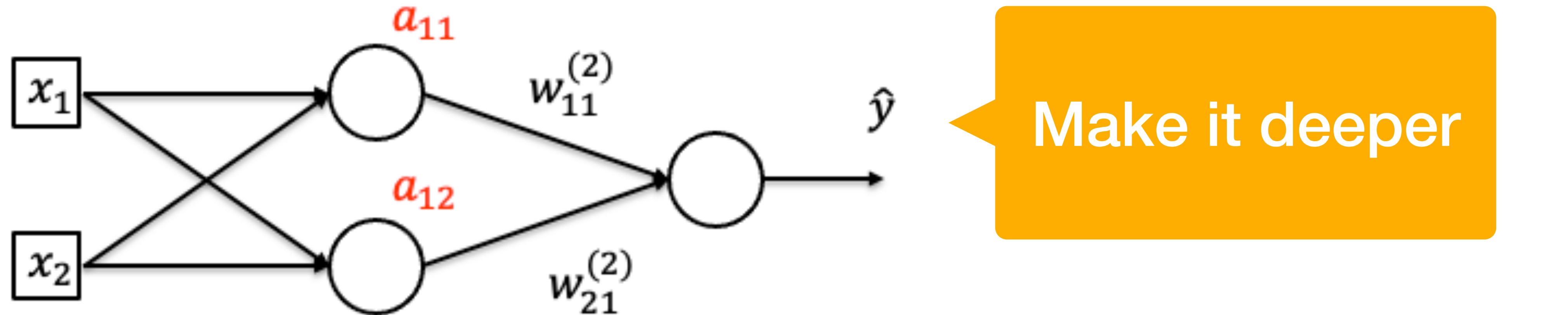
Below the graph, the gradient of the loss function with respect to the input  $x_1$  is derived using the chain rule:

$$\frac{\partial l}{\partial x_1} = \frac{\partial l}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} w_{11} = (\hat{y} - y)w_{11}$$

- By chain rule:

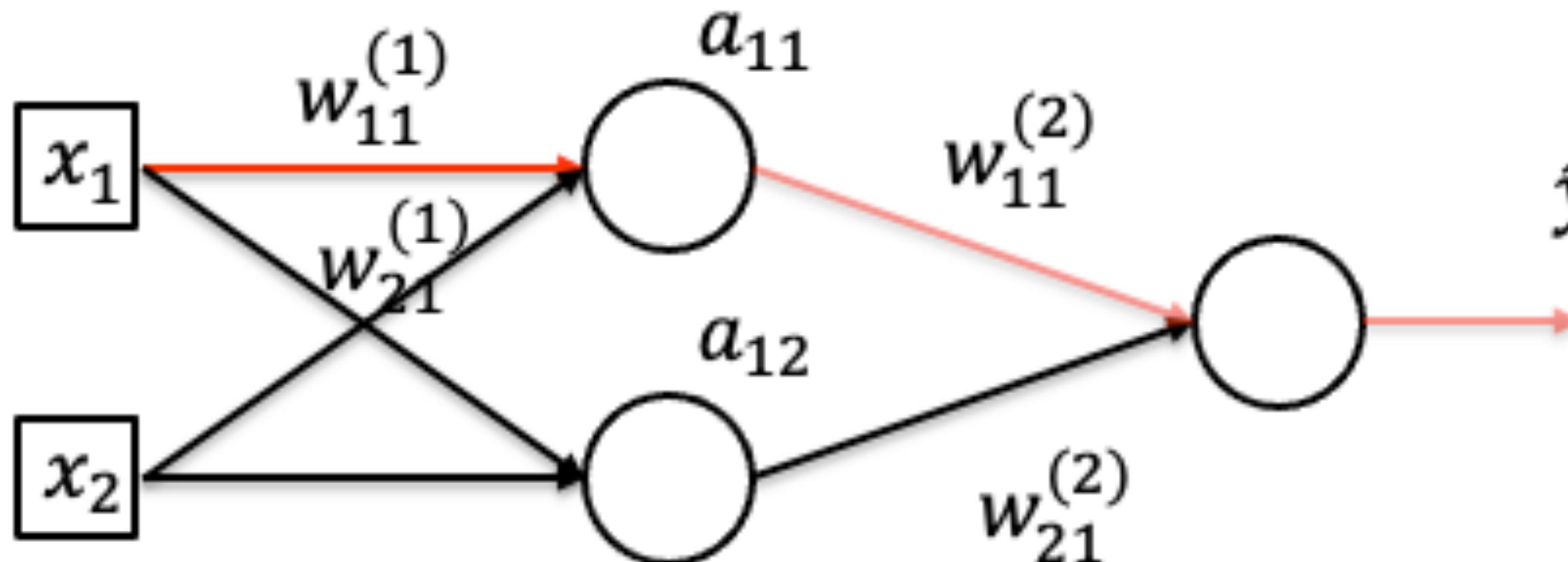
$$\frac{\partial l}{\partial x_1} = \frac{\partial l}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} w_{11} = (\hat{y} - y)w_{11}$$

# Calculate Gradient (on one data point)



- By chain rule:  $\frac{\partial l}{\partial a_{11}} = (\hat{y} - y)w_{11}^{(2)}$ ,  $\frac{\partial l}{\partial a_{12}} = (\hat{y} - y)w_{21}^{(2)}$

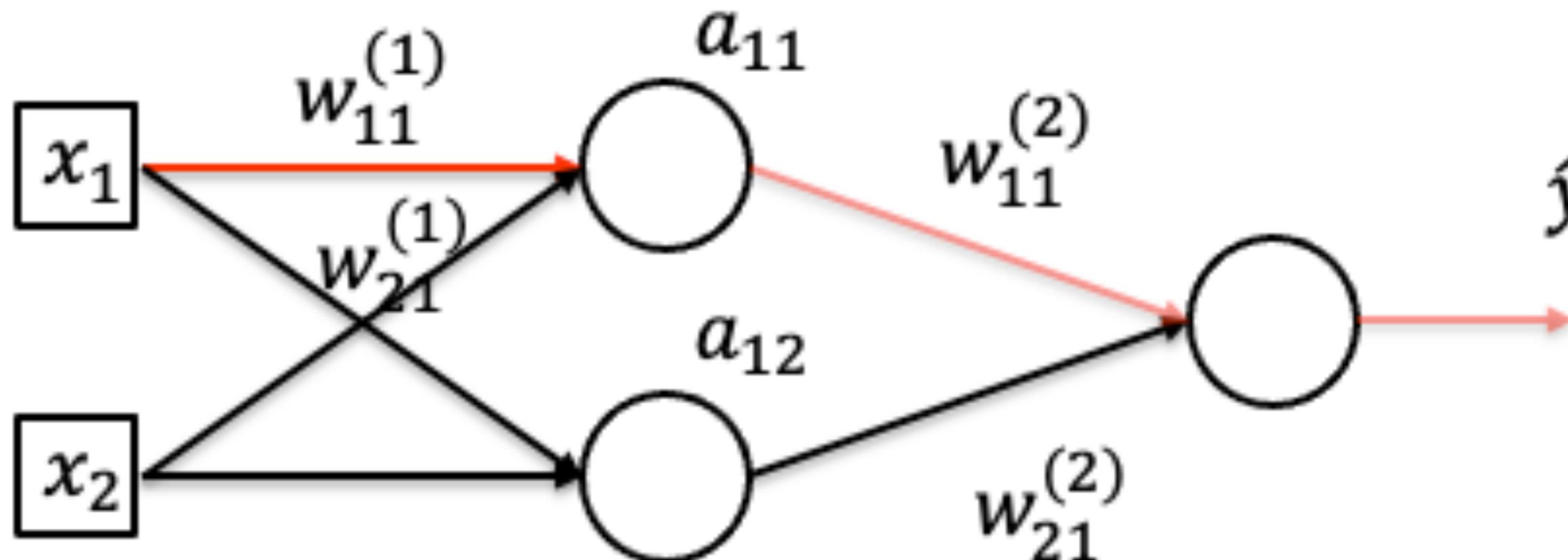
# Calculate Gradient (on one data point)



$$\begin{aligned} & w_{11}^{(1)} x_1 \\ & w_{21}^{(1)} x_2 \end{aligned} \rightarrow \begin{array}{c} + \\ \text{---} \end{array} \rightarrow z_{11} \xrightarrow{\sigma(z_{11})} a_{11} \xrightarrow{\frac{\partial l}{\partial a_{11}}} l(x, y)$$
$$\frac{\partial a_{11}}{\partial z_{11}} = \sigma'(z_{11}) \quad \frac{\partial l}{\partial a_{11}} = (\hat{y} - y) w_{11}^{(2)}$$

- By chain rule:  $\frac{\partial l}{\partial w_{11}^{(1)}} = \frac{\partial l}{\partial a_{11}} \frac{\partial a_{11}}{\partial w_{11}^{(1)}} = (\hat{y} - y) w_{11}^{(2)} \frac{\partial a_{11}}{\partial w_{11}^{(1)}}$

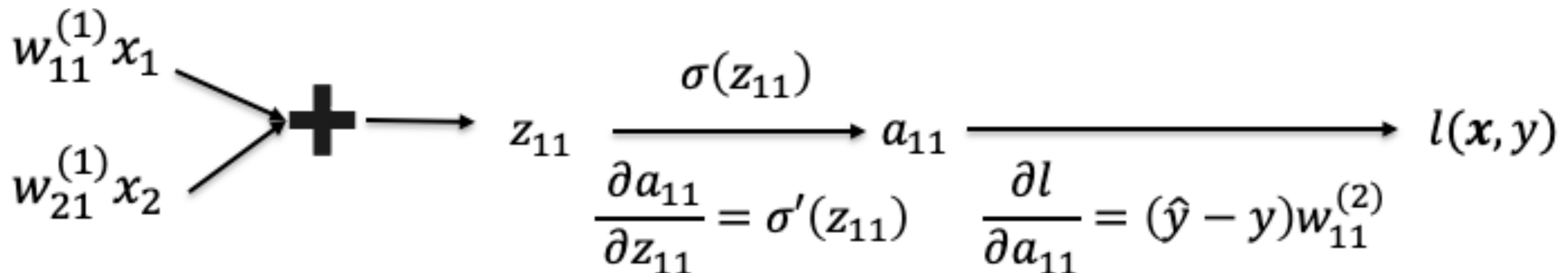
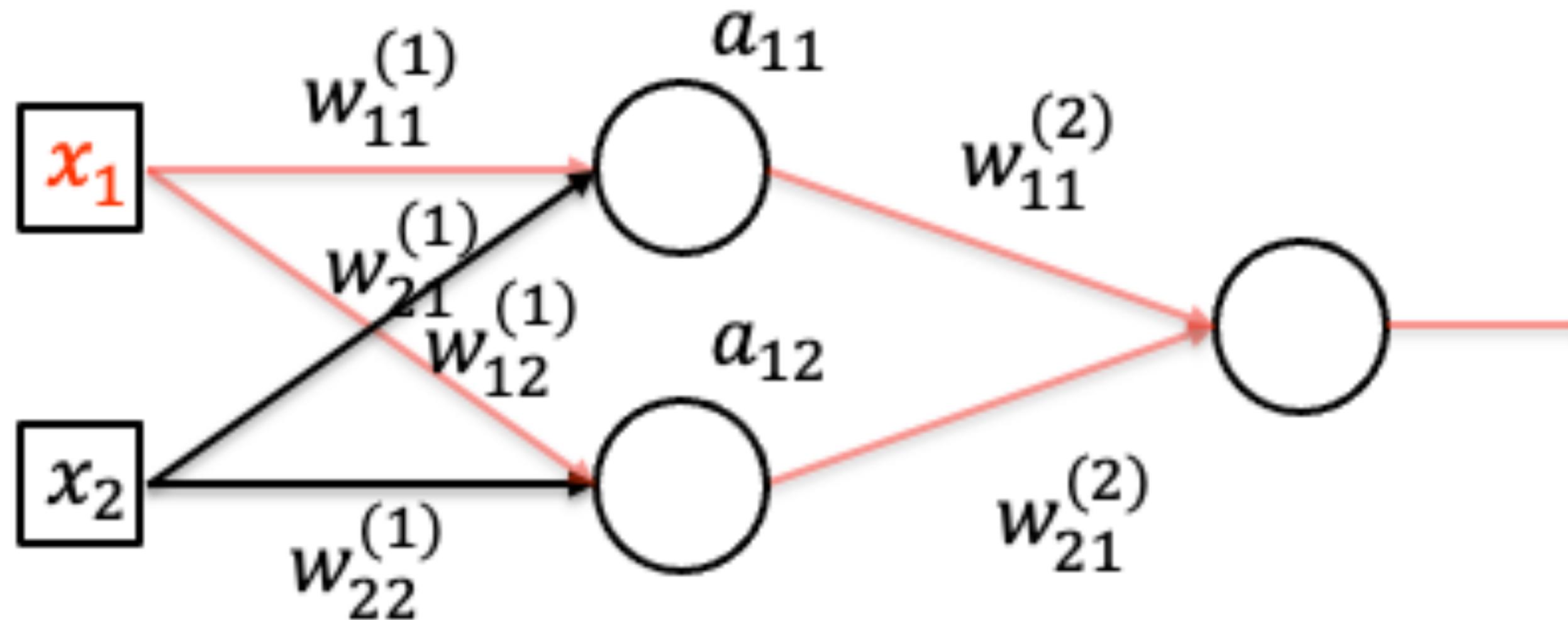
# Calculate Gradient (on one data point)



$$\begin{aligned} & w_{11}^{(1)} x_1 \quad \quad \quad z_{11} \quad \quad \quad a_{11} \quad \quad \quad l(x, y) \\ & w_{21}^{(1)} x_2 \quad \quad \quad \sigma(z_{11}) \quad \quad \quad \frac{\partial a_{11}}{\partial z_{11}} = \sigma'(z_{11}) \quad \quad \quad \frac{\partial l}{\partial a_{11}} = (\hat{y} - y) w_{11}^{(2)} \end{aligned}$$

- By chain rule:  $\frac{\partial l}{\partial w_{11}^{(1)}} = \frac{\partial l}{\partial a_{11}} \frac{\partial a_{11}}{\partial w_{11}^{(1)}} = (\hat{y} - y) w_{11}^{(2)} a_{11} (1 - a_{11}) x_1$

# Calculate Gradient (on one data point)



- By chain rule:

$$\frac{\partial l}{\partial x_1} = \frac{\partial l}{\partial a_{11}} \frac{\partial a_{11}}{\partial x_1} + \frac{\partial l}{\partial a_{12}} \frac{\partial a_{12}}{\partial x_1}$$

# Quiz Break

Gradient Descent in neural network training computes the \_\_\_\_\_ of a loss function with respect to the model \_\_\_\_\_ until convergence.

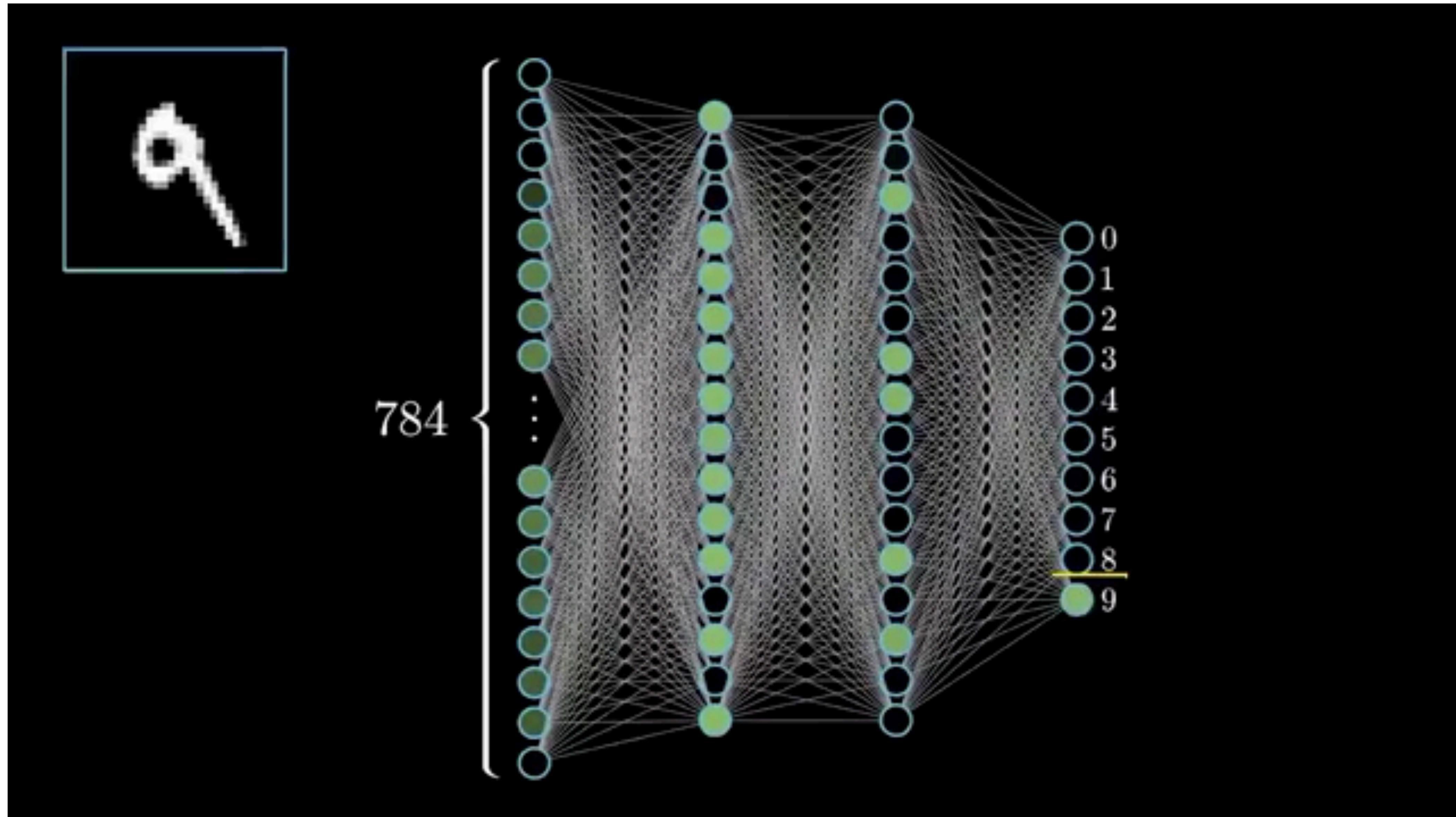
- A gradients, parameters
- B parameters, gradients
- C loss, parameters
- D parameters, loss

# Quiz Break

Suppose you are given a dataset with 1,000,000 images to train with. Which of the following methods is more desirable if training resources are limit but decent accuracy is needed?

- A Gradient Descent
- B Stochastic Gradient Descent
- C Minibatch Stochastic Gradient Descent
- D Computation Graph

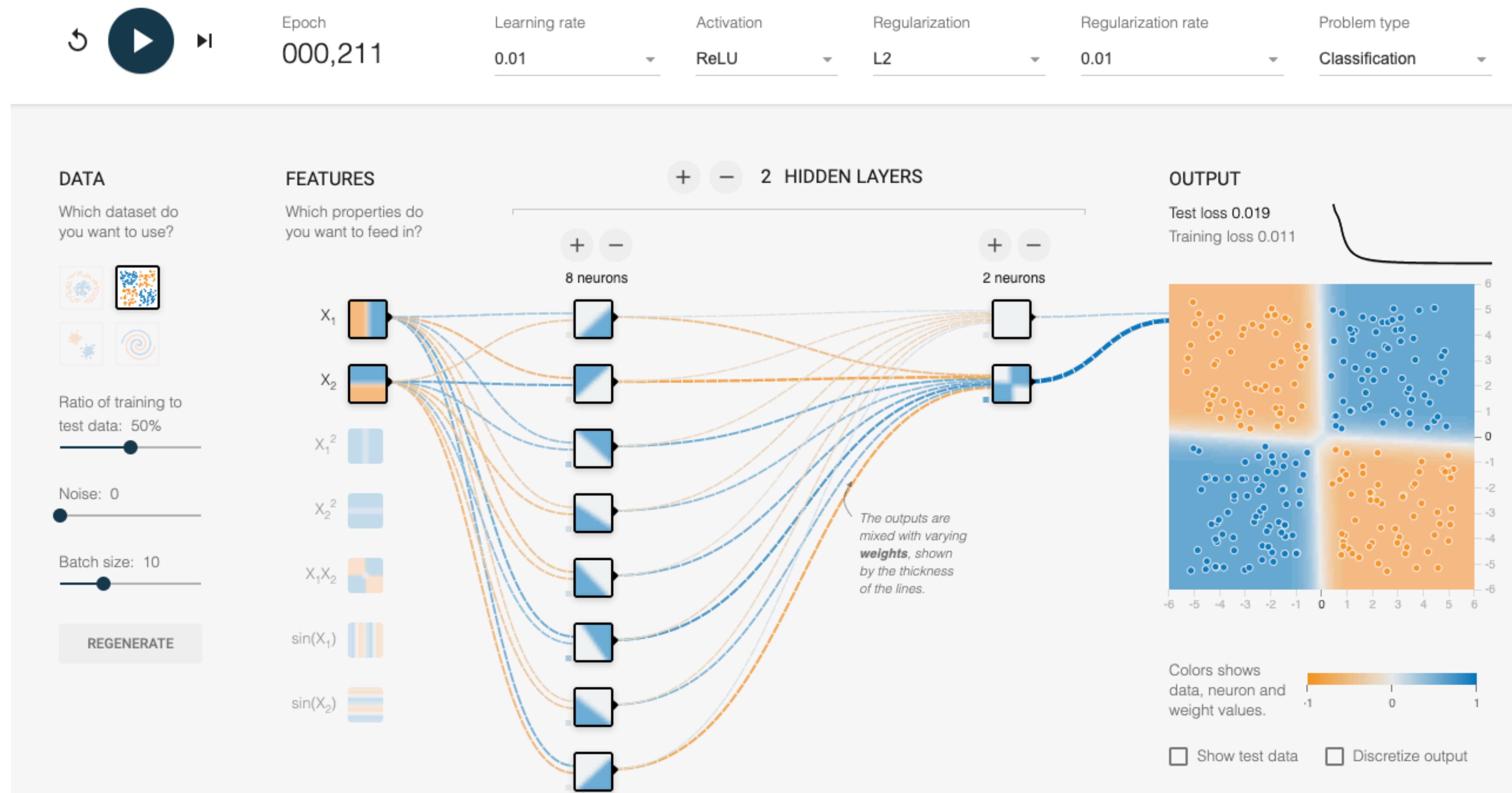
# HW6



# HW6 (working with MNIST dataset)



# Demo: Learning XOR using neural net



• <https://playground.tensorflow.org/>

# What we've learned today...

- Single-layer Perceptron Review
- Multi-layer Perceptron
  - Single output
  - Multiple output
- How to train neural networks
  - Gradient descent



**Thanks!**