# CS 540 Introduction to Artificial Intelligence
## Neural Networks (II)
University of Wisconsin-Madison

**Fall 2022**

# Today's outline

- Single-layer Perceptron Review

- 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:

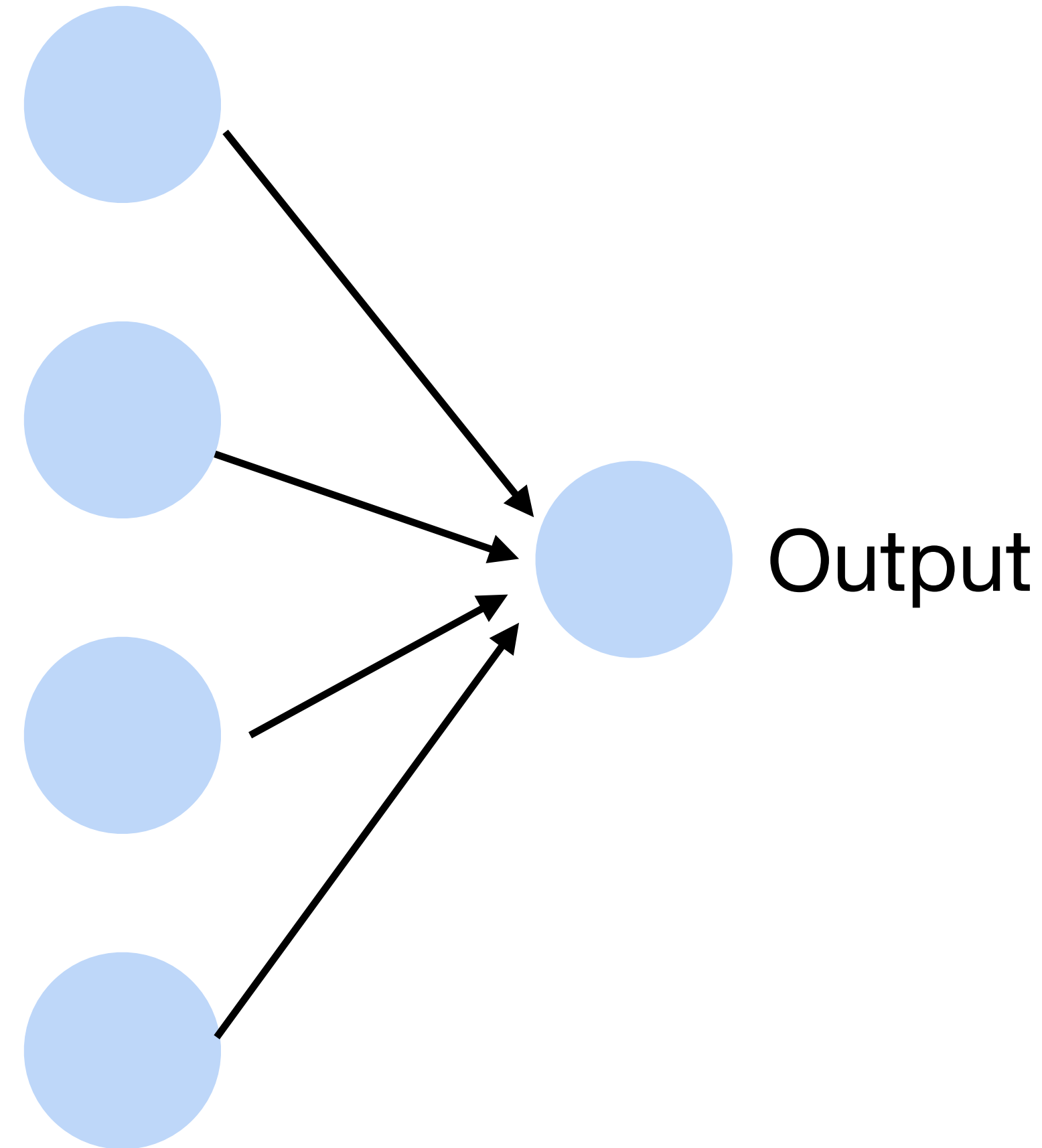$$o = \sigma\left(\langle \mathbf{w}, \mathbf{x}\rangle + b\right)$$

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

**Activation function**
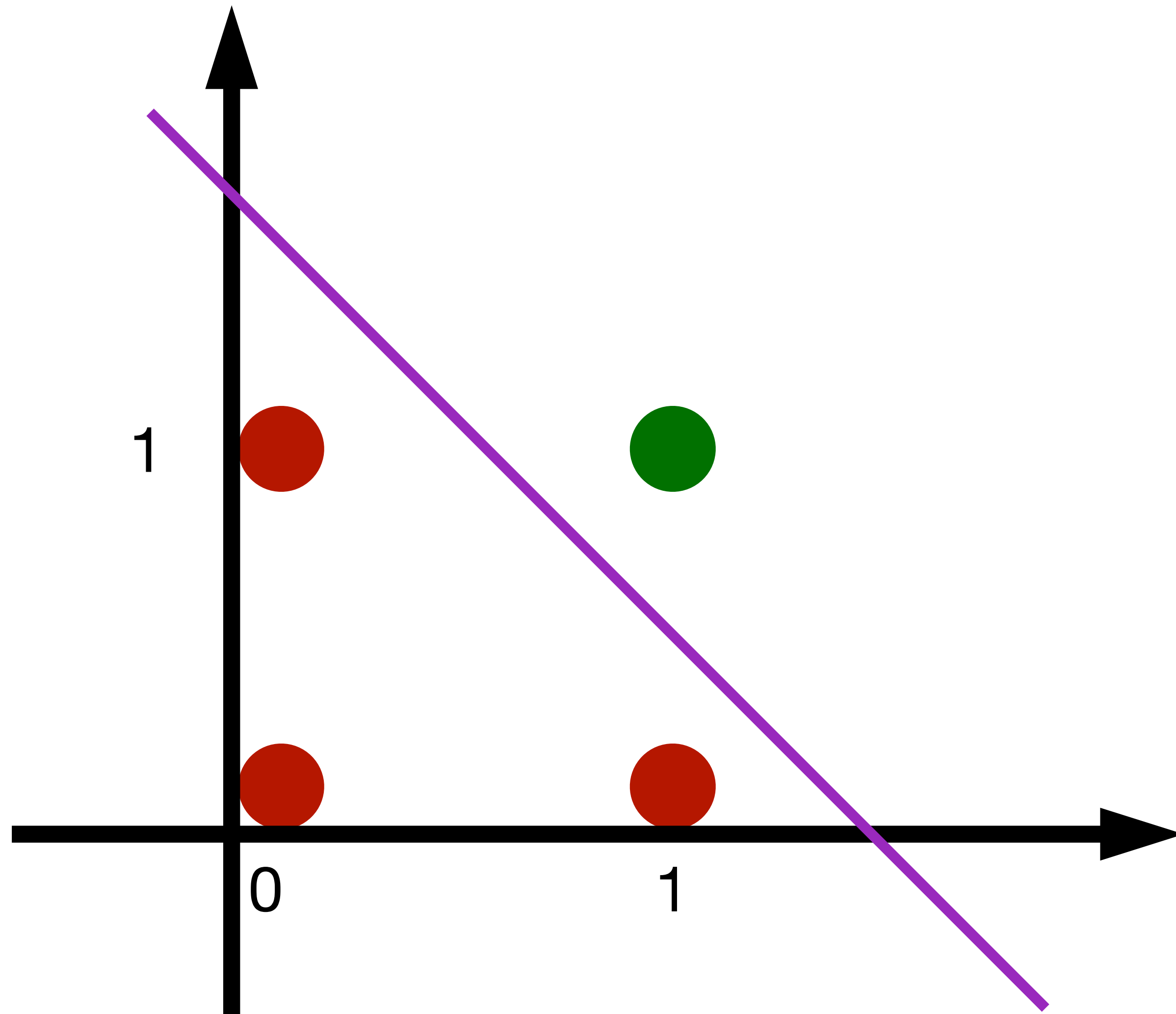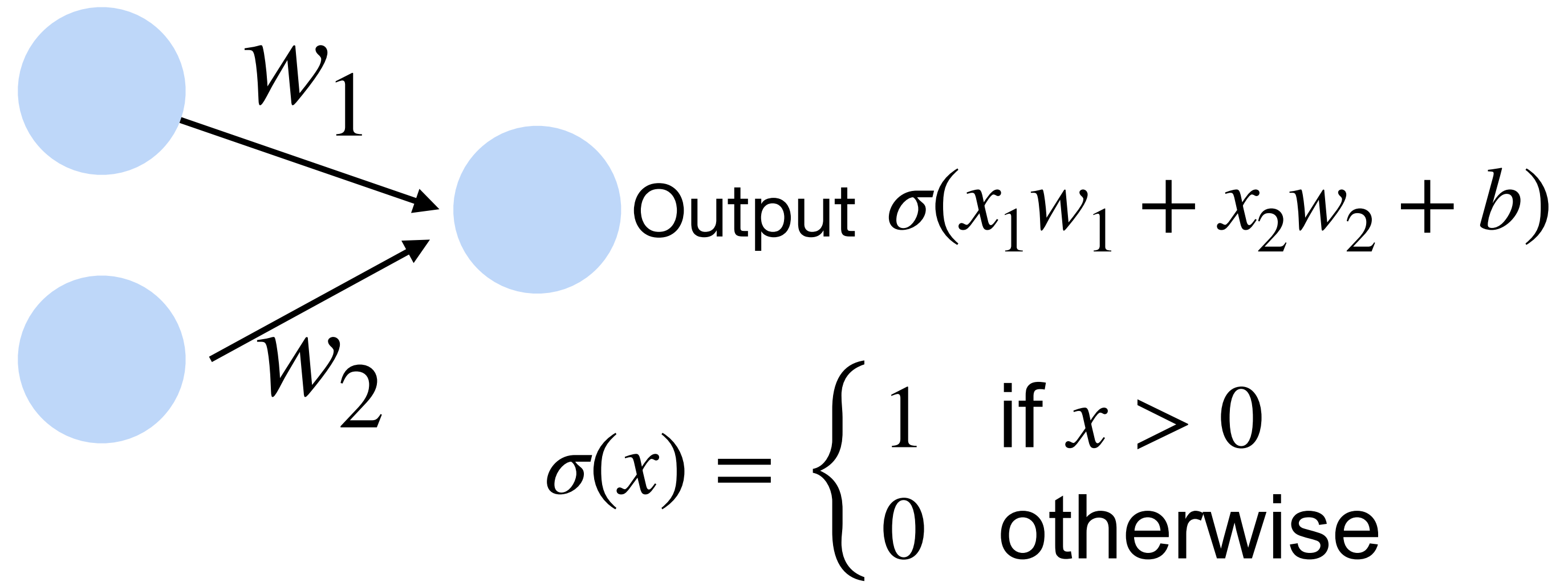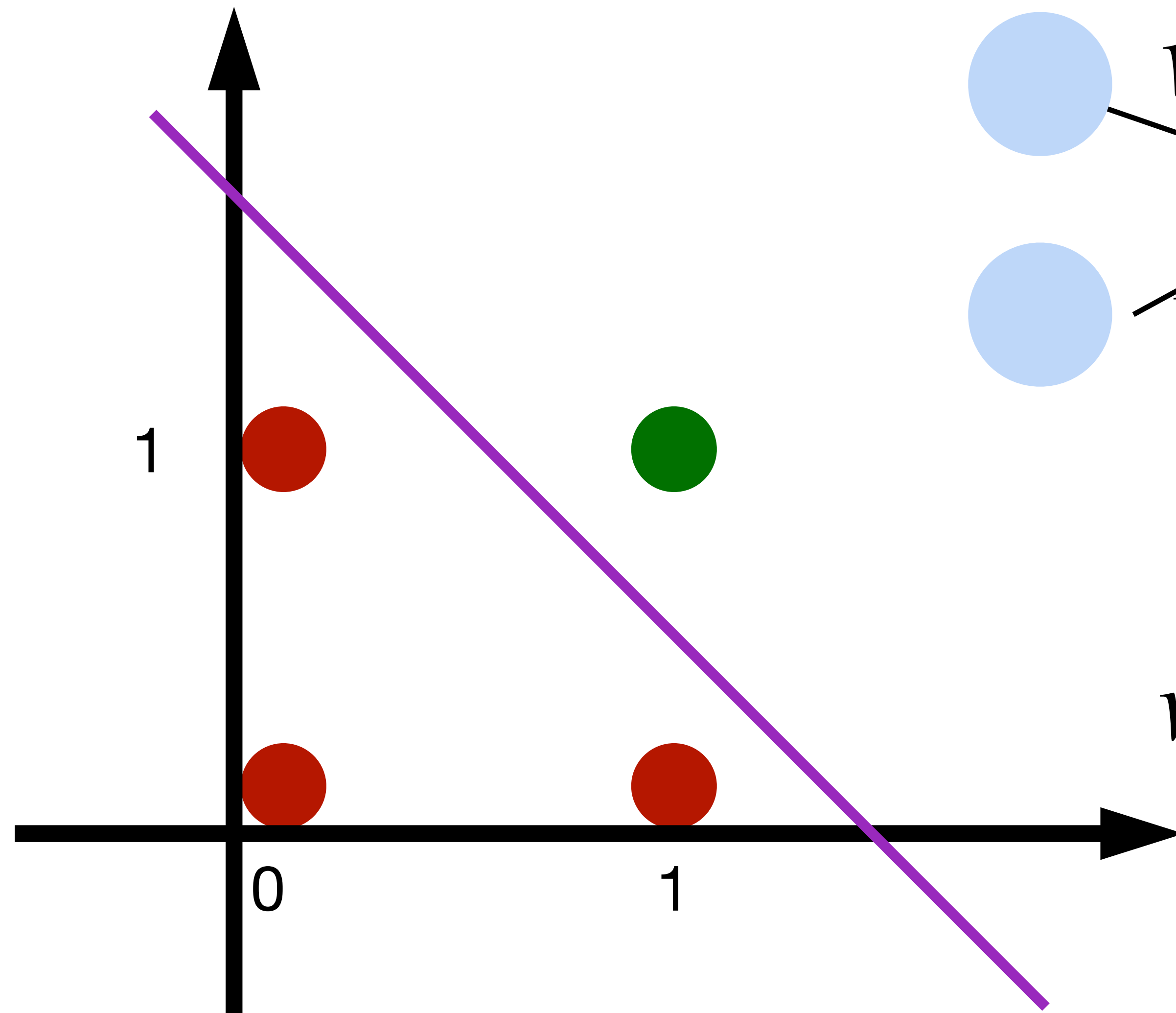
**Cats vs. dogs?**

Input

Output

# Learning AND function using perceptron

The perceptron can learn an AND function

# Learning AND function using perceptron

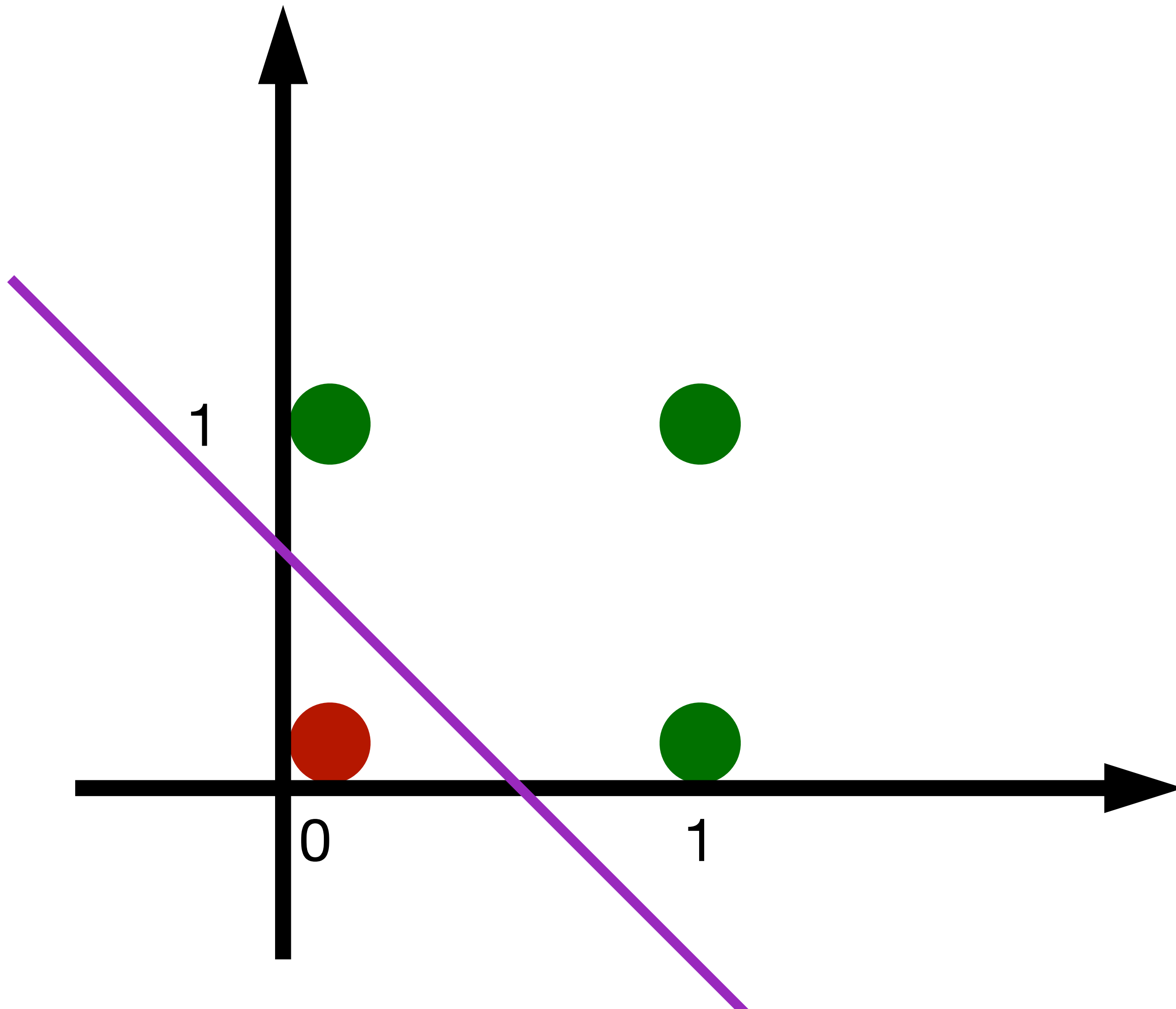The perceptron can learn an AND function



Output $\sigma(x_1 w_1 + x_2 w_2 + b)$

$$\sigma(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$$
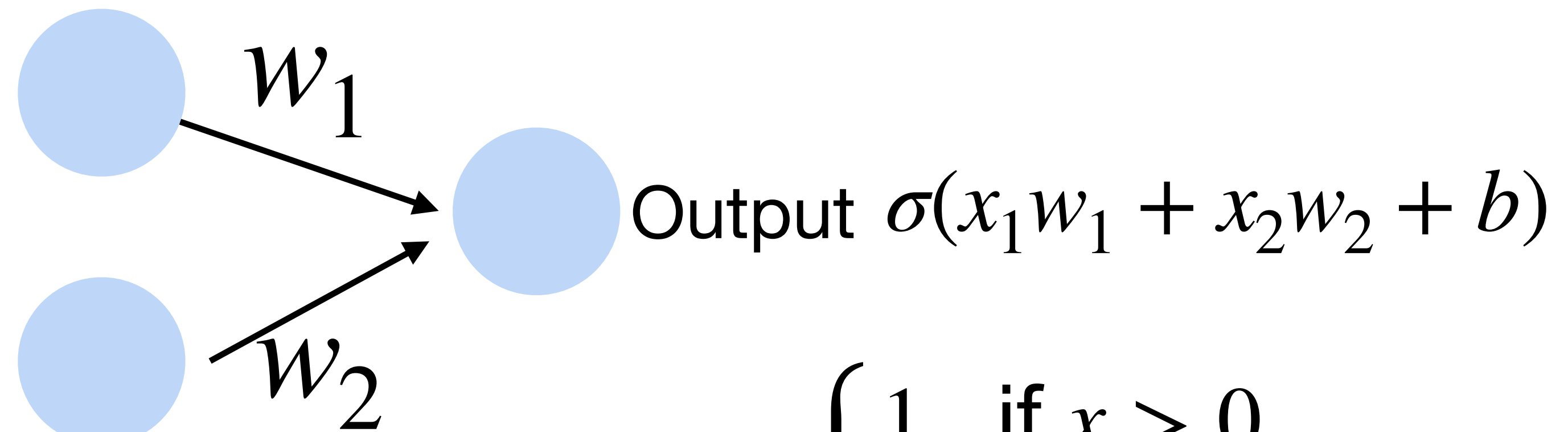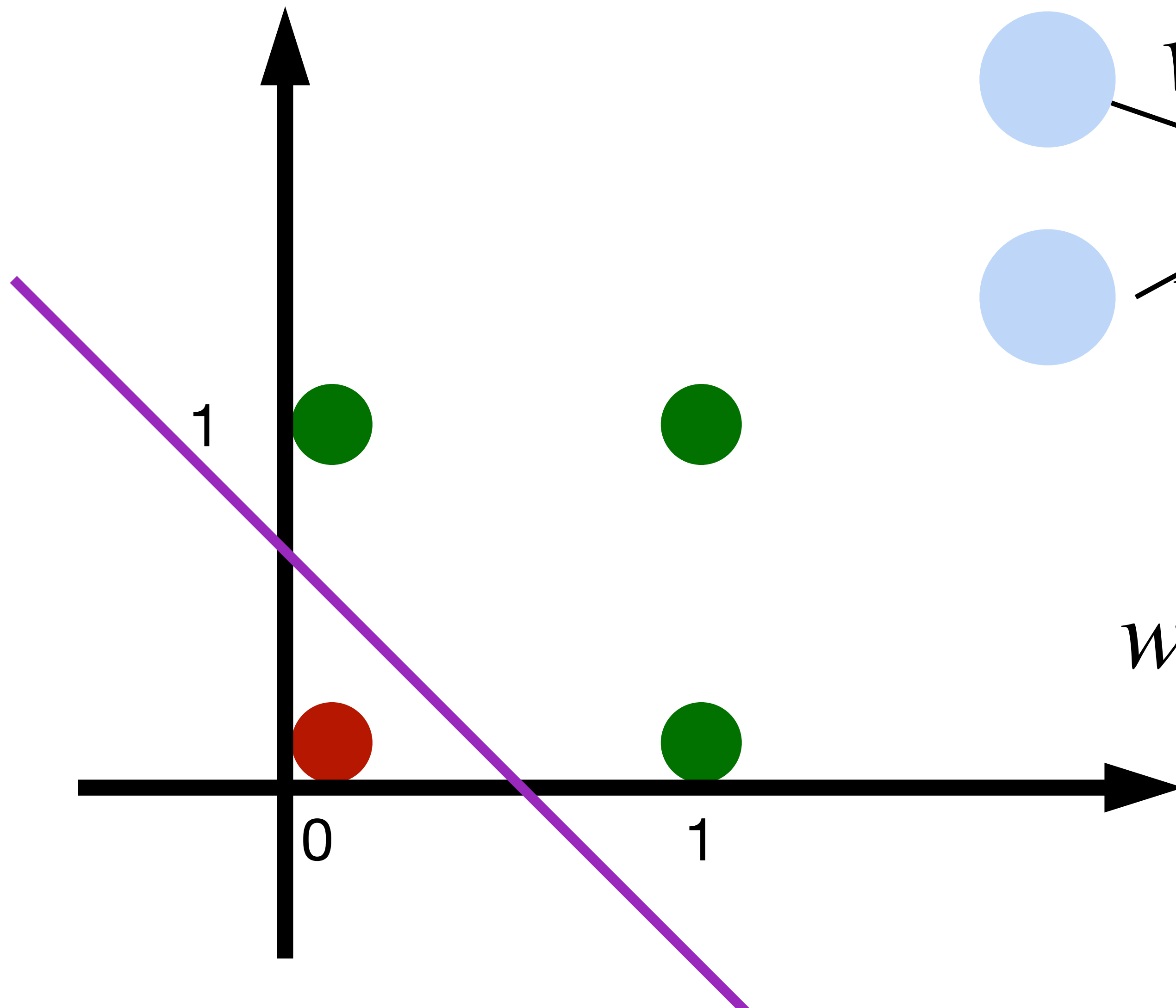
$$w_1 = 1, w_2 = 1, b = -1.5$$

# Learning OR function using perceptron

The perceptron can learn an OR function

# Learning OR function using perceptron

The perceptron can learn an OR function



Output $\sigma(x_1 w_1 + x_2 w_2 + b)$

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

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

# Learning NOT function using perceptron

The perceptron can learn NOT function (single input)

$$x \xrightarrow{w_1} \quad \text{Output } \sigma(xw_1 + b)$$

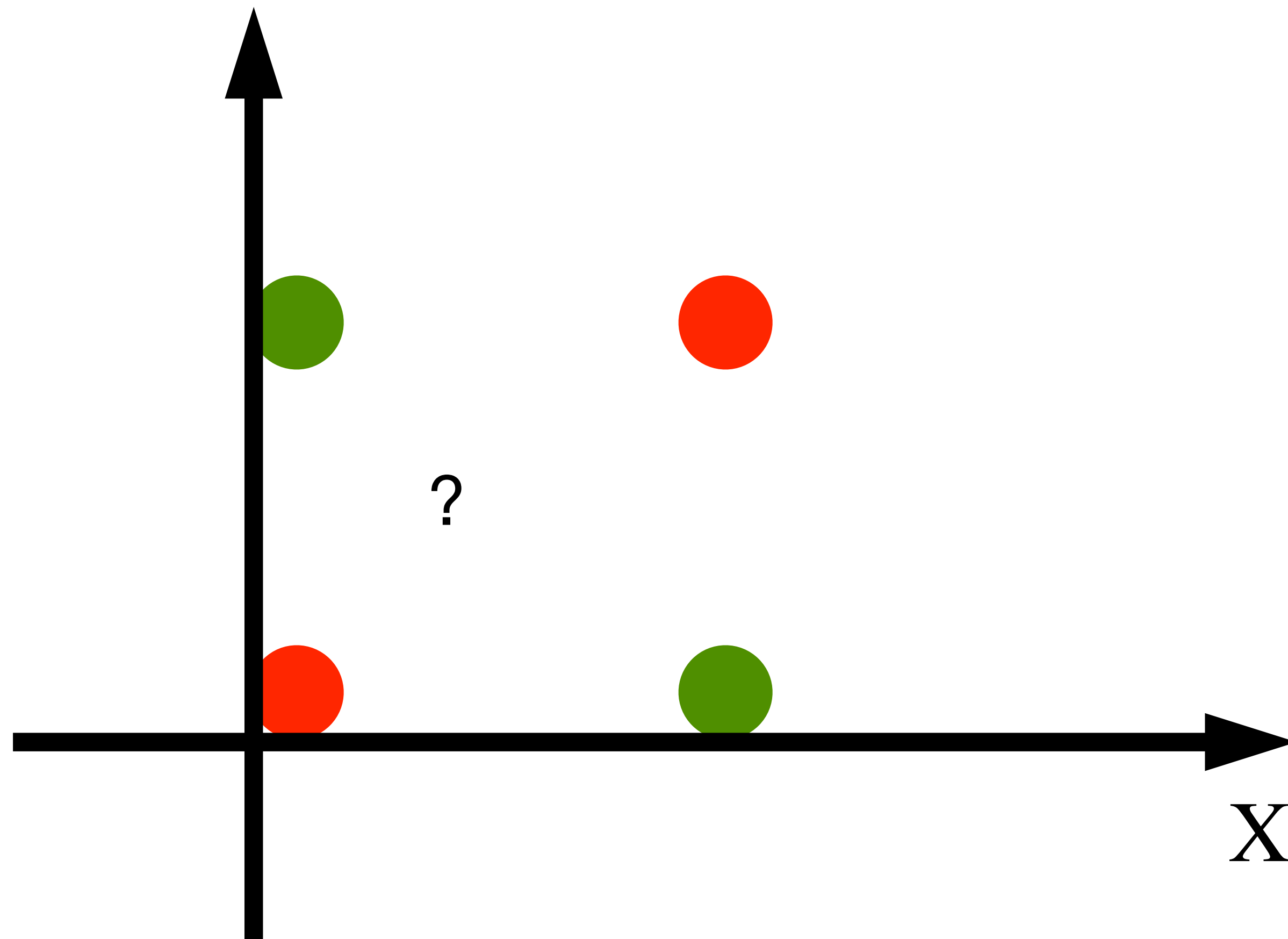$$\sigma(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$$

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

# The limited power of a single neuron

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



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

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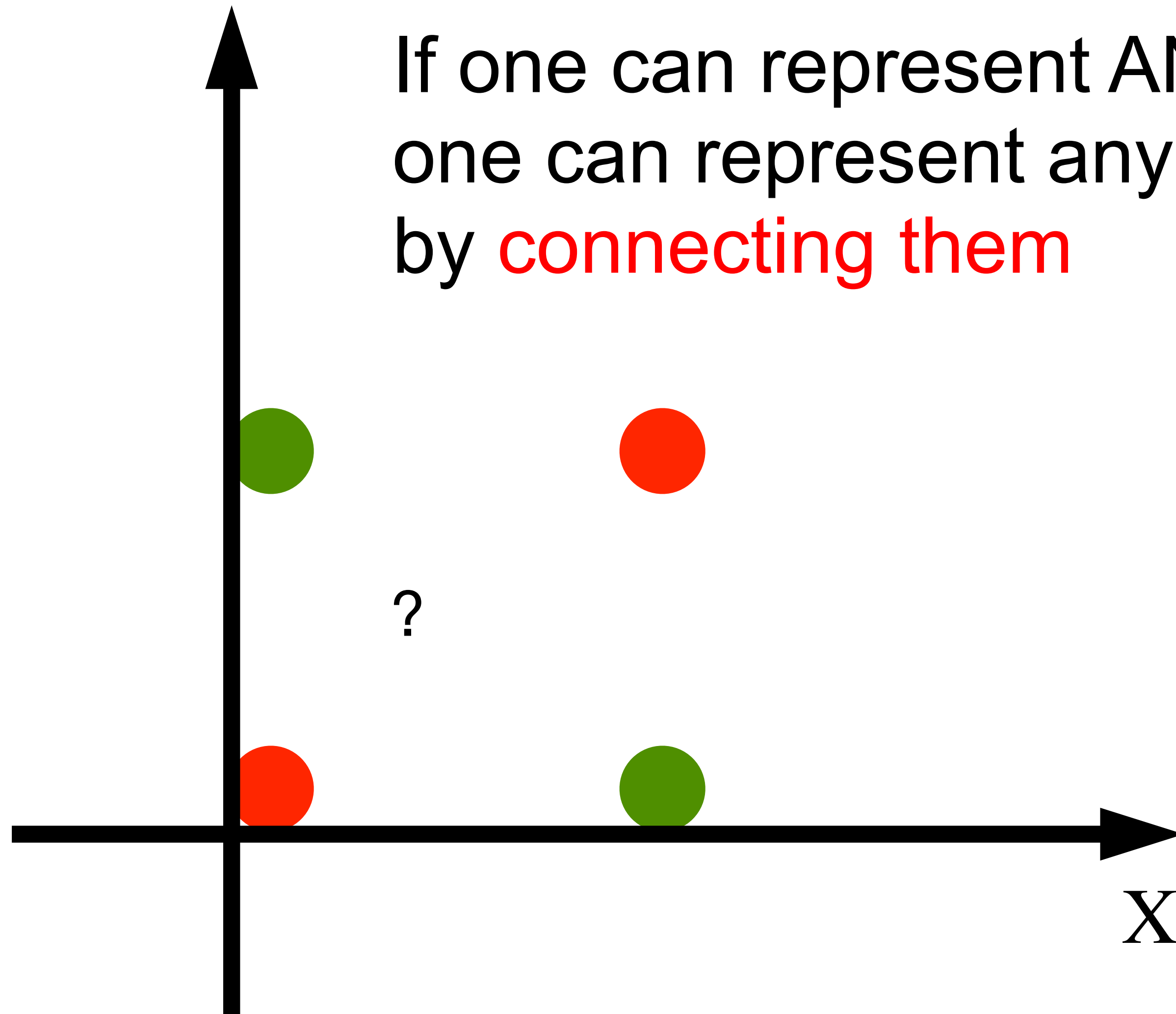
$$x_1 = 0, x_2 = 1, y = 1$$

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

$$XOR(x1, x2) = (x1 \land \neg x2) \lor (\neg x1 \land x2)$$

# The limited power of a single neuron

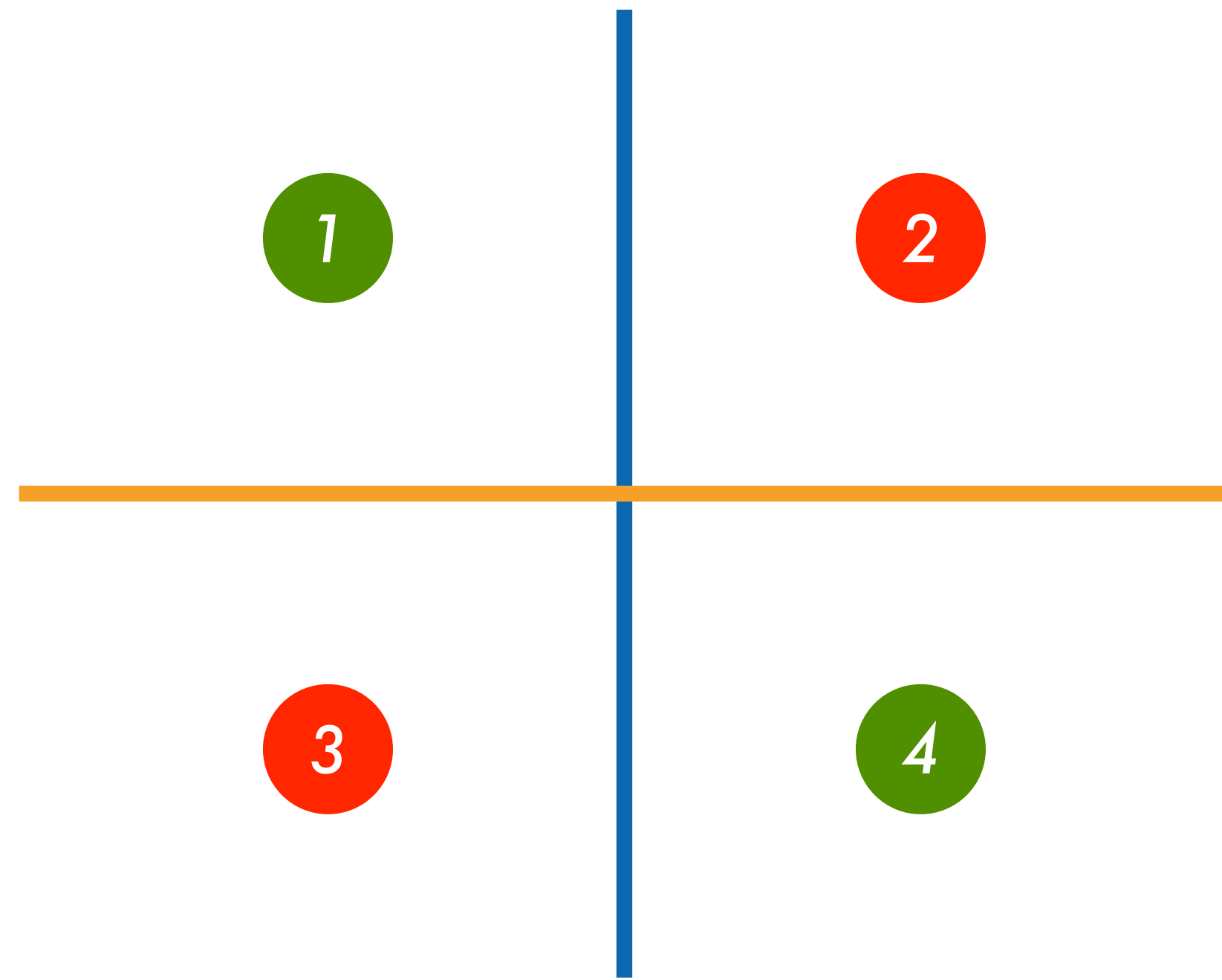**XOR** problem

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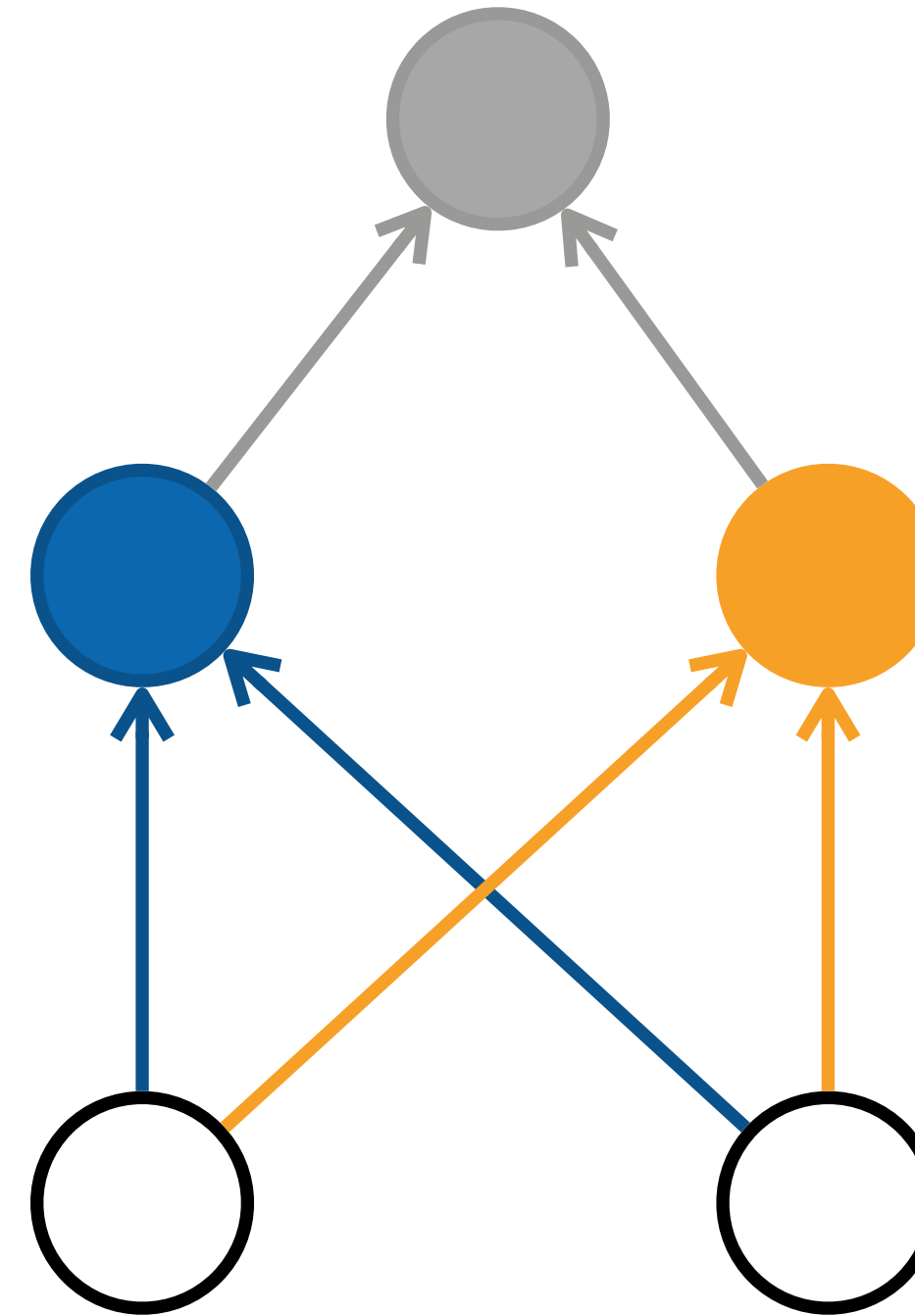one can represent any logic circuit (including XOR),
by connecting them

?

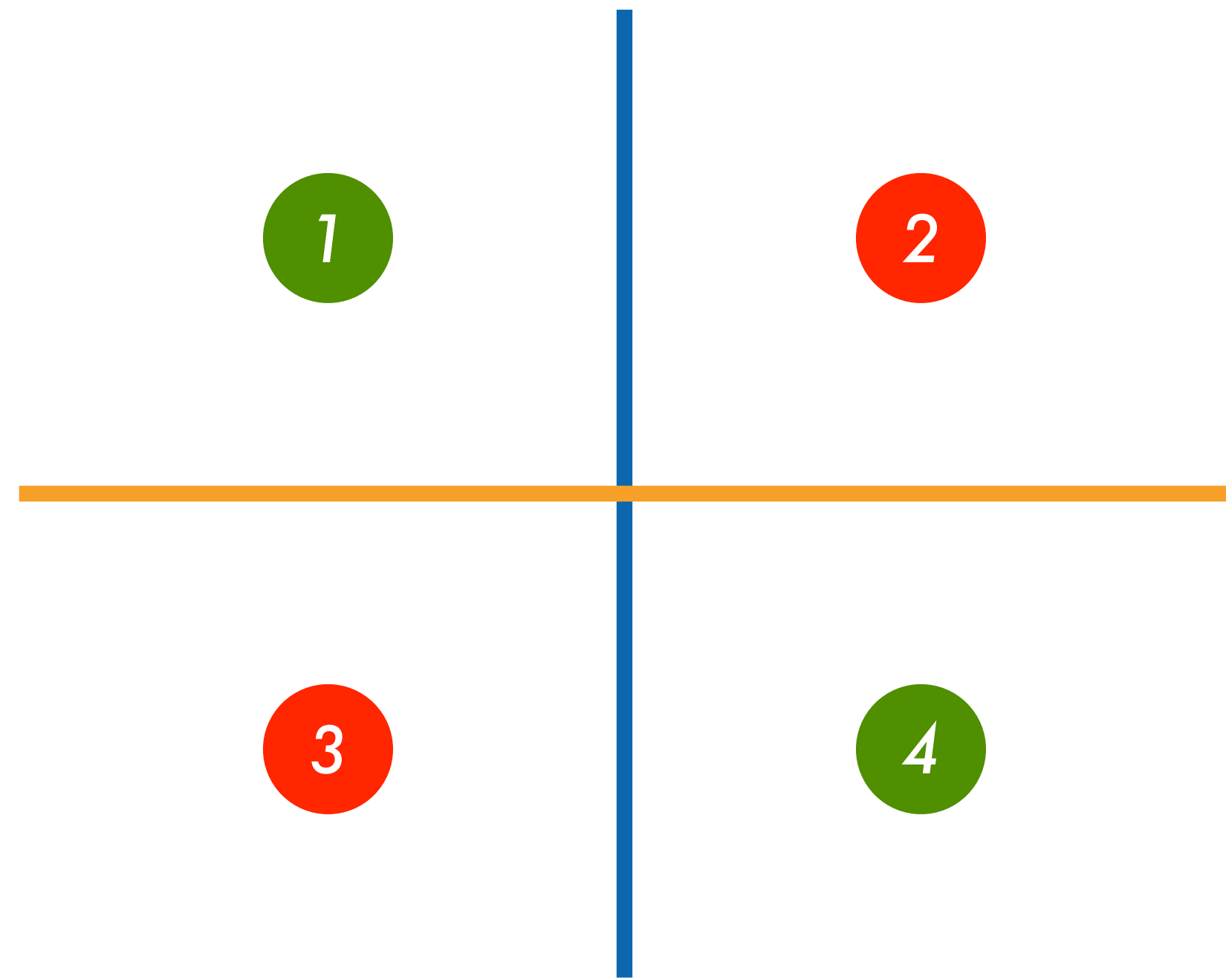$$XOR(x1,x2) = (x1 \wedge \neg x2) \vee (\neg x1 \wedge x2)$$
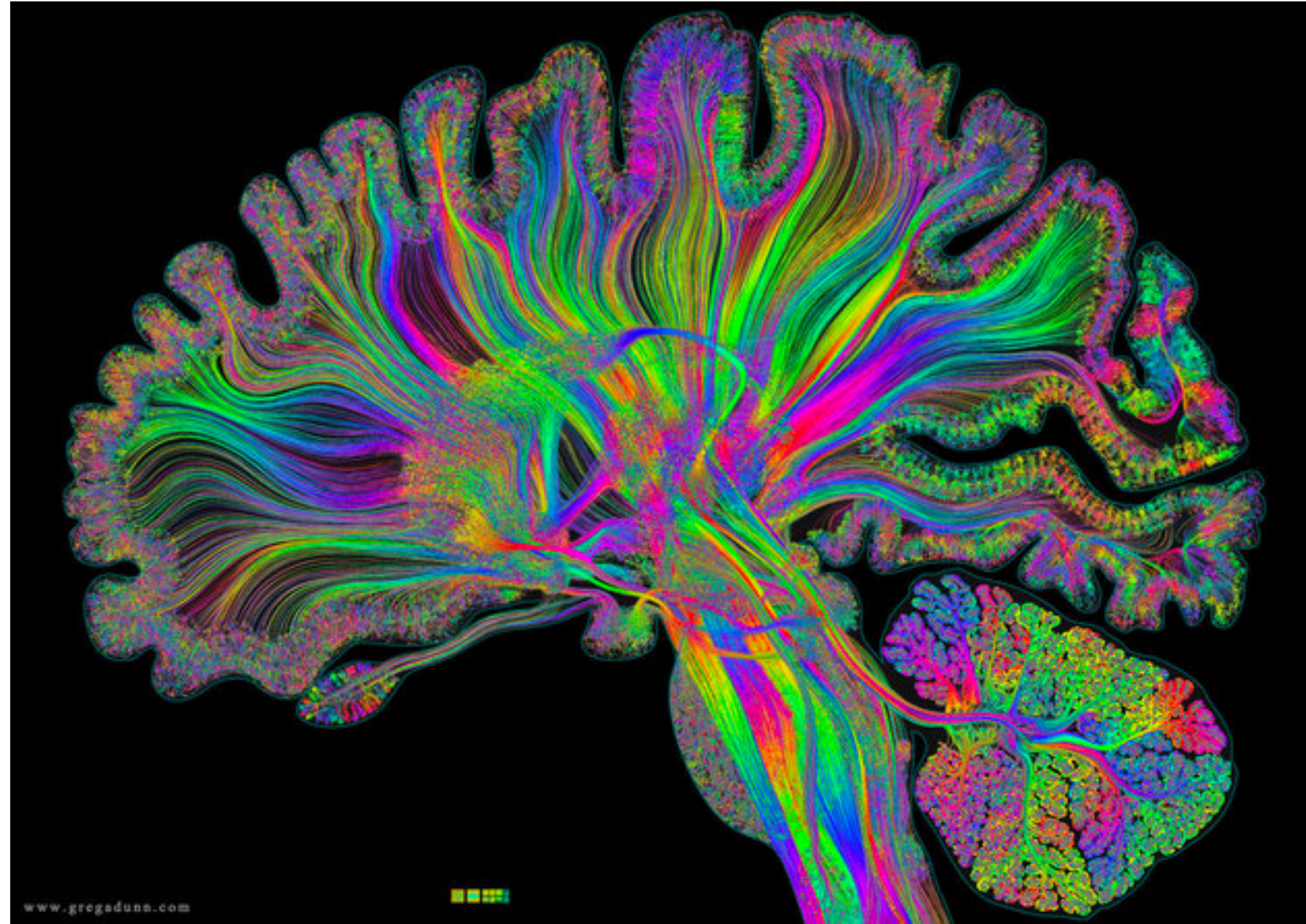
# Learning XOR

# Learning XOR

# Multilayer Perceptron

# Multi-layer perceptron: Example

- Standard way to connect Perceptrons
- Example: 1 hidden layer, 1 output layer, depth = 2

Input

Hidden layer
3 neurons

$$h_1 = \sigma(\sum_{i=1}^{d} x_i w_{1i}^{(1)} + b_1)$$

$\mathbf{x} \in \mathbb{R}^d$

$w_{11}^{(1)}$

$w_{12}^{(1)}$

$x_1$

$x_2$

# Multi-layer perceptron: Example

- Standard way to connect Perceptrons
- Example: 1 hidden layer, 1 output layer, depth = 2

Hidden layer
3 neurons

Input

$\mathbf{x} \in \mathbb{R}^d$

$x_1$

$w_{21}^{(1)}$

$w_{22}^{(1)}$

$x_2$

$$h_2 = \sigma\left(\sum_{i=1}^{d} x_i w_{2i}^{(1)} + b_2\right)$$
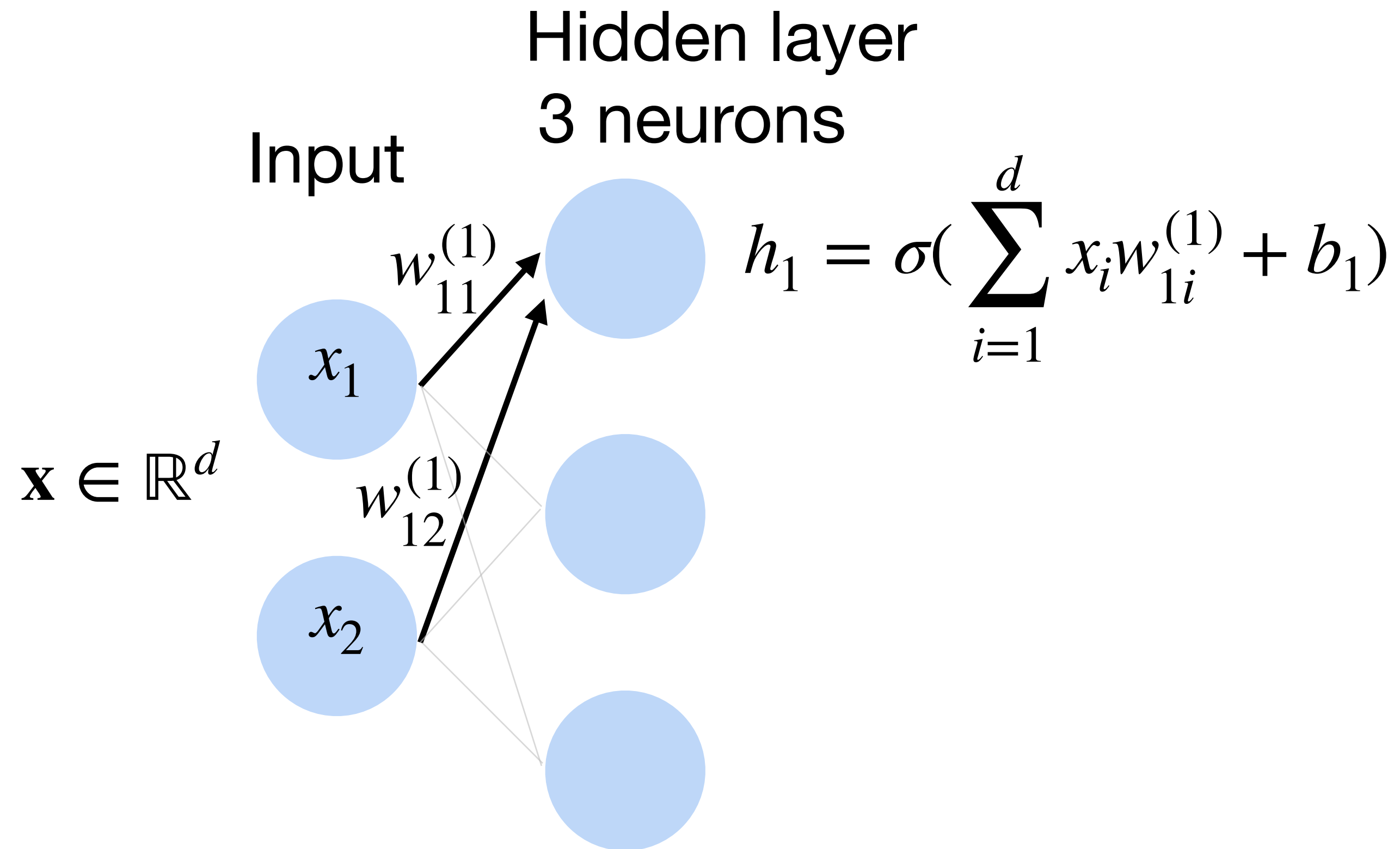
# Multi-layer perceptron: Example

- Standard way to connect Perceptrons
- Example: 1 hidden layer, 1 output layer, depth = 2

Hidden layer
3 neurons

Input

$x_1$

$x_2$

$\mathbf{x} \in \mathbb{R}^d$

$w_{31}^{(1)}$

$w_{32}^{(1)}$

$$h_3 = \sigma(\sum_{i=1}^{d} x_i w_{3i}^{(1)} + b_3)$$

# Multi-layer perceptron: Example

- Standard way to connect Perceptrons
- Example: 1 hidden layer, 1 output layer, depth = 2

Input

Hidden layer
m=3 neurons

$\mathbf{x} \in \mathbb{R}^d$

$x_1$

$x_2$

$h_1 = \sigma(\sum_{i=1}^{d} x_i w_{1i}^{(1)} + b_1)$

$h_2 = \sigma(\sum_{i=1}^{d} x_i w_{2i}^{(1)} + b_2)$
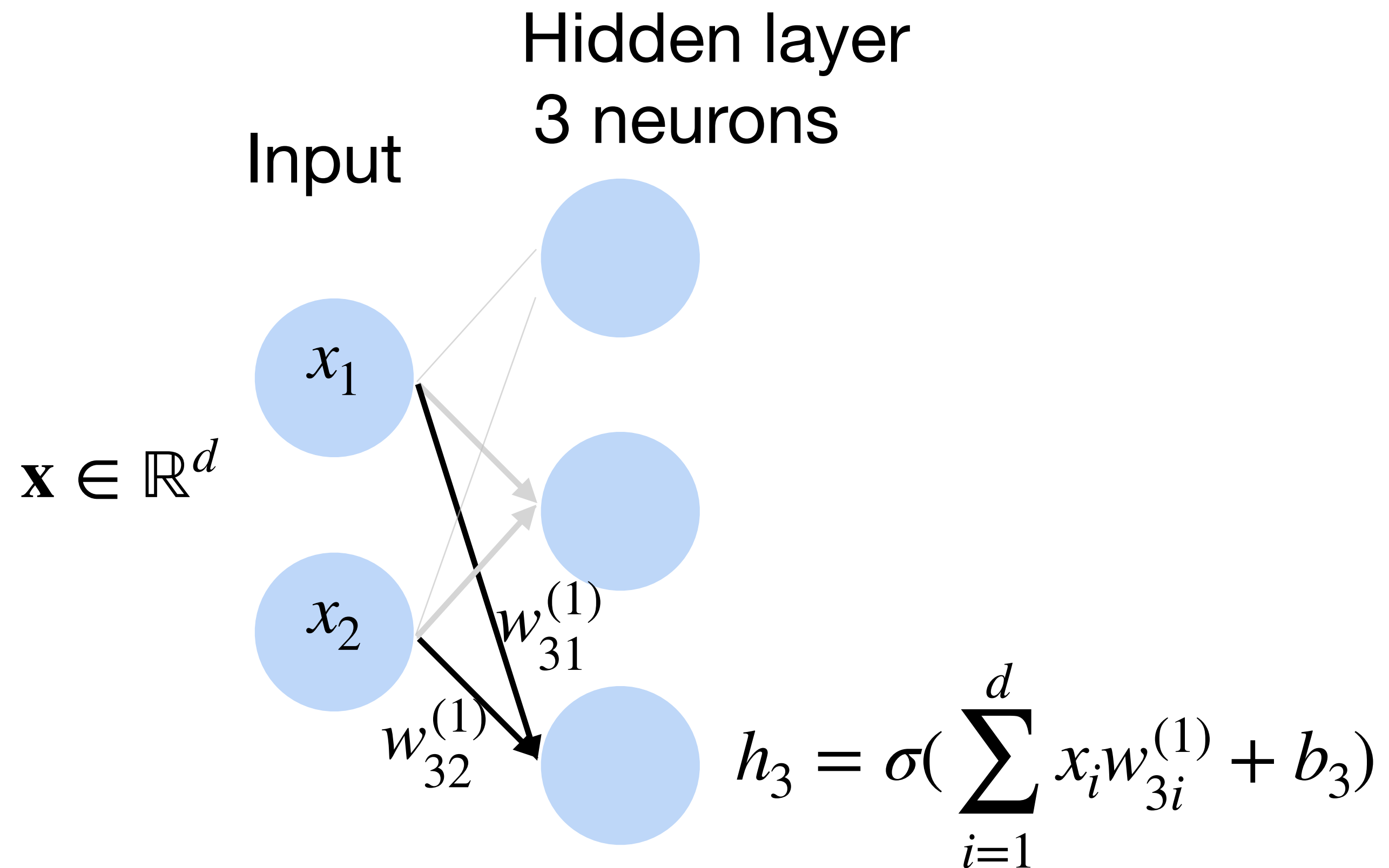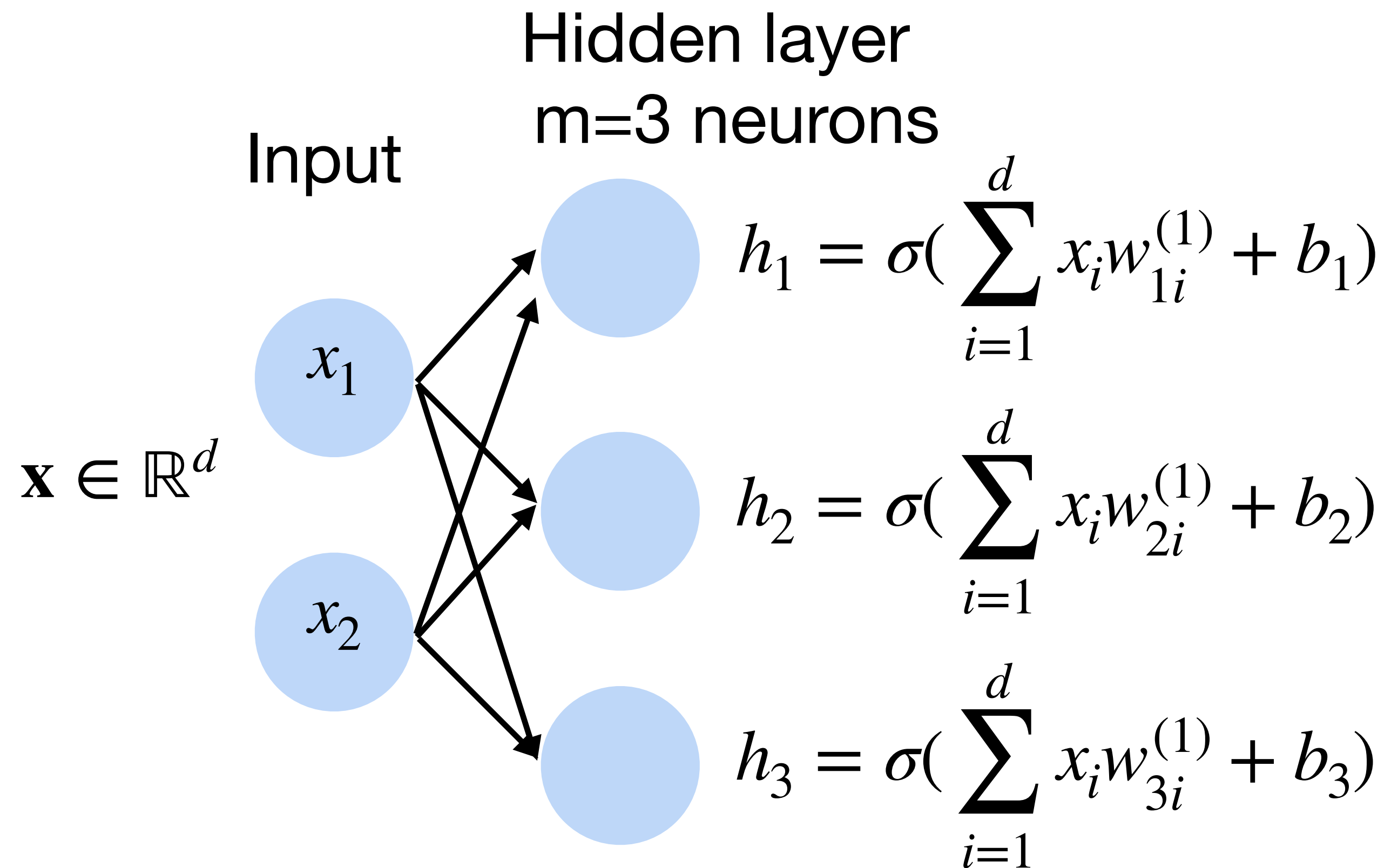
$h_3 = \sigma(\sum_{i=1}^{d} x_i w_{3i}^{(1)} + b_3)$

# Multi-layer perceptron: Example

- Standard way to connect Perceptrons
- Example: 1 hidden layer, 1 output layer, depth = 2

Input

Hidden layer
m=3 neurons

$\mathbf{x} \in \mathbb{R}^d$

$x_1$

$x_2$

$h_1 = \sigma(\sum_{i=1}^{d} x_i w_{1i}^{(1)} + b_1)$

$h_2 = \sigma(\sum_{i=1}^{d} x_i w_{2i}^{(1)} + b_2)$

$h_3 = \sigma(\sum_{i=1}^{d} x_i w_{3i}^{(1)} + b_3)$

$w_1^{(2)}$

$w_2^{(2)}$

$w_3^{(2)}$

Output

# Multi-layer perceptron: Example

- Standard way to connect Perceptrons
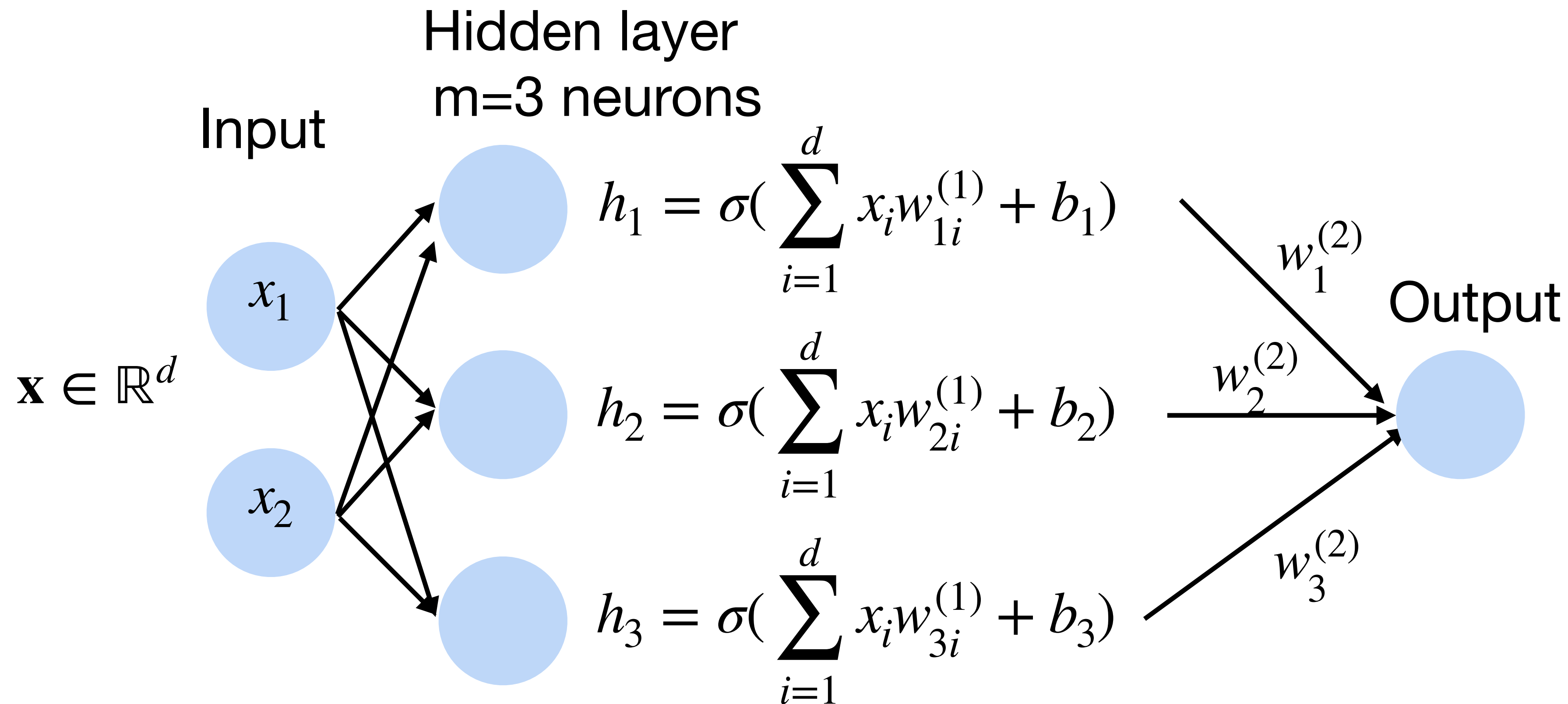- Example: 1 hidden layer, 1 output layer, depth = 2

Input

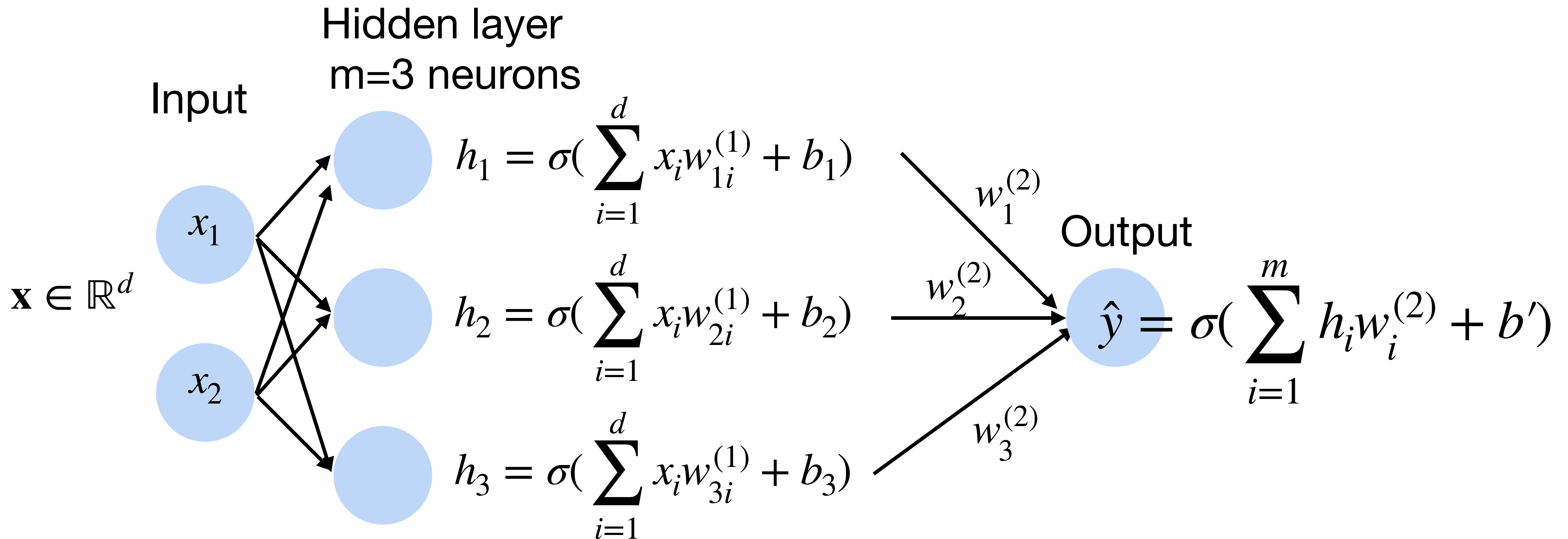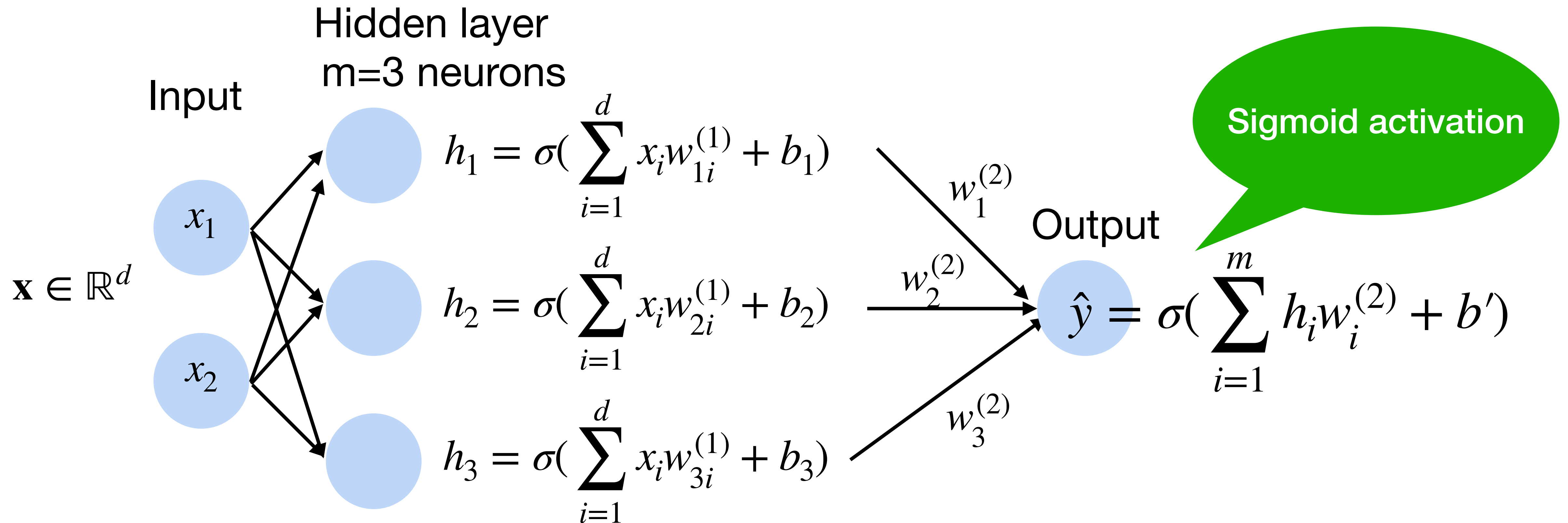Hidden layer
m=3 neurons

Output

$$\mathbf{x} \in \mathbb{R}^d$$

$x_1$

$x_2$

$h_1 = \sigma(\sum_{i=1}^{d} x_i w_{1i}^{(1)} + b_1)$

$h_2 = \sigma(\sum_{i=1}^{d} x_i w_{2i}^{(1)} + b_2)$

$h_3 = \sigma(\sum_{i=1}^{d} x_i w_{3i}^{(1)} + b_3)$

$w_1^{(2)}$

$w_2^{(2)}$

$w_3^{(2)}$

$$\hat{y} = \sigma(\sum_{i=1}^{m} h_i w_i^{(2)} + b')$$
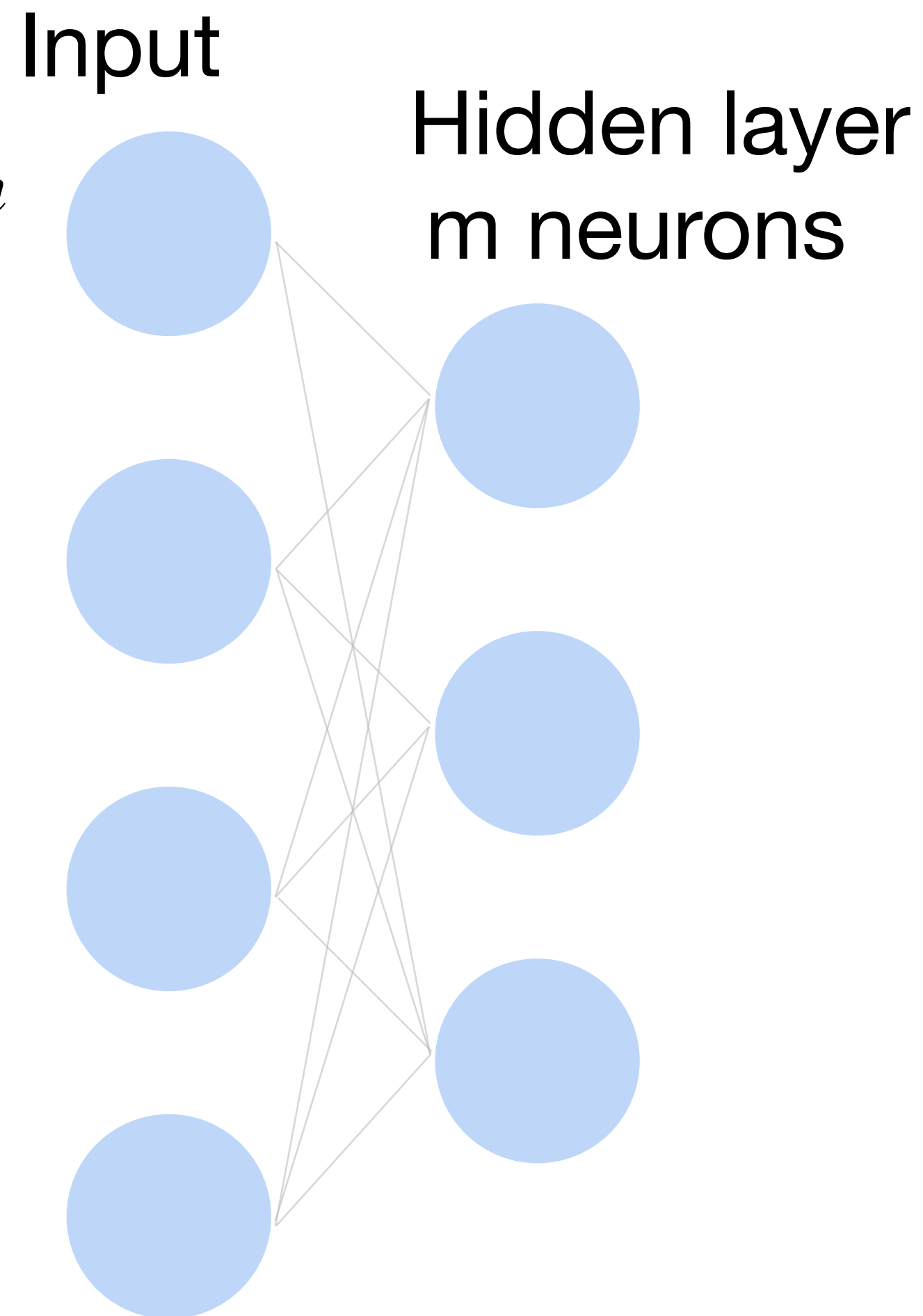
# Multi-layer perceptron: Example

- Standard way to connect Perceptrons
- Example: 1 hidden layer, 1 output layer, depth = 2



Input

Hidden layer
m=3 neurons

$\mathbf{x} \in \mathbb{R}^d$

$x_1$

$x_2$

$h_1 = \sigma(\sum_{i=1}^{d} x_i w_{1i}^{(1)} + b_1)$

$h_2 = \sigma(\sum_{i=1}^{d} x_i w_{2i}^{(1)} + b_2)$

$h_3 = \sigma(\sum_{i=1}^{d} x_i w_{3i}^{(1)} + b_3)$

$w_1^{(2)}$

$w_2^{(2)}$

$w_3^{(2)}$

Output

Sigmoid activation

$\hat{y} = \sigma(\sum_{i=1}^{m} h_i w_i^{(2)} + b')$
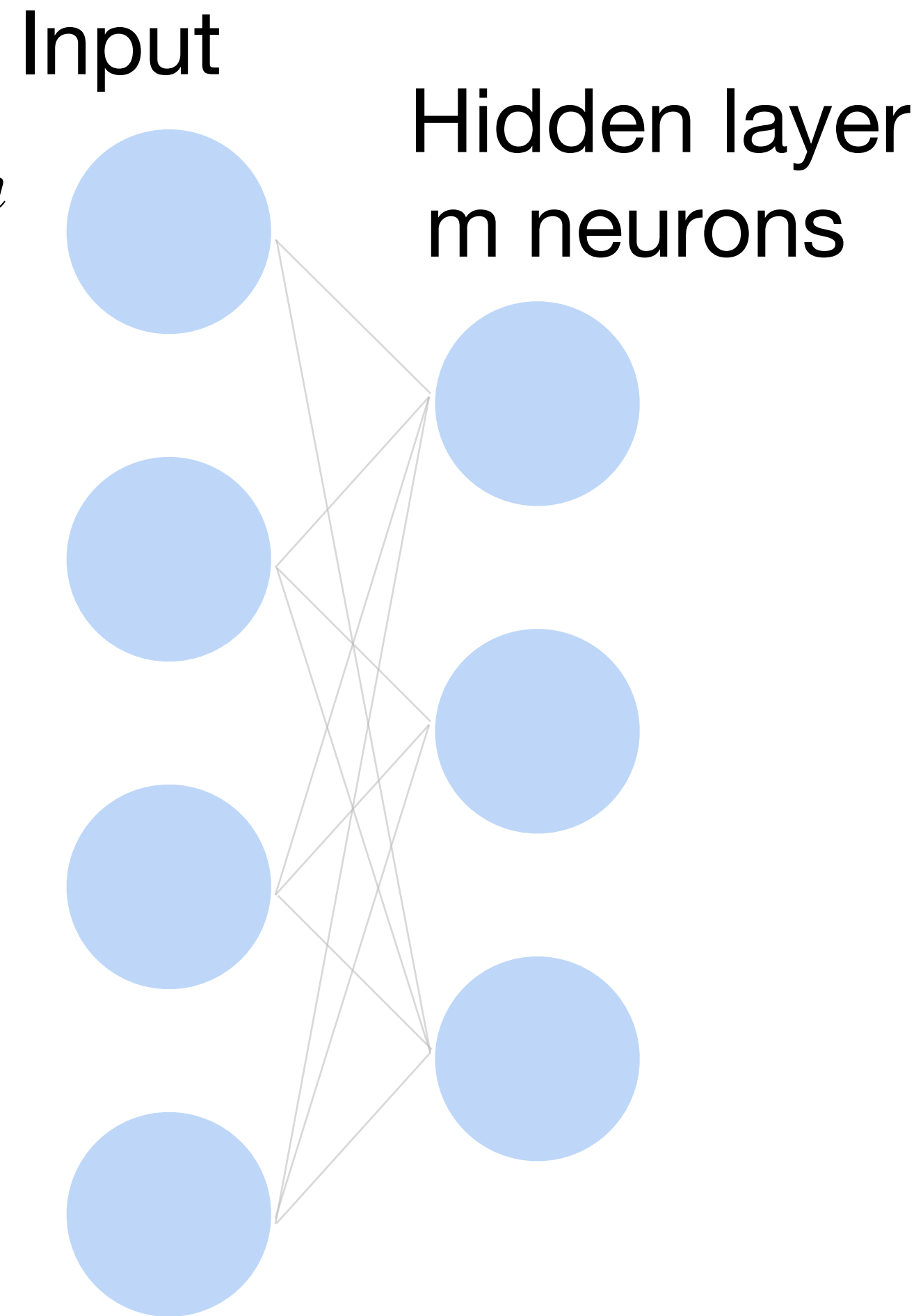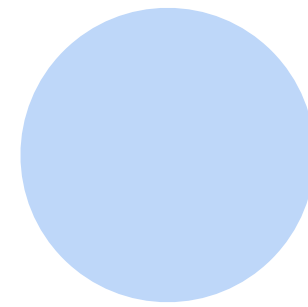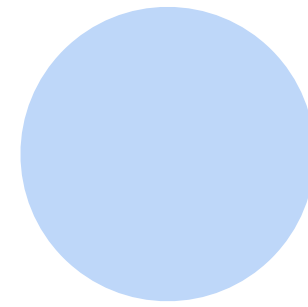
# Multi-layer perceptron: Matrix Notation

- Input $\mathbf{x} \in \mathbb{R}^d$
- Hidden $\mathbf{W}^{(1)} \in \mathbb{R}^{m \times d}, \mathbf{b} \in \mathbb{R}^m$
- Intermediate output

Input

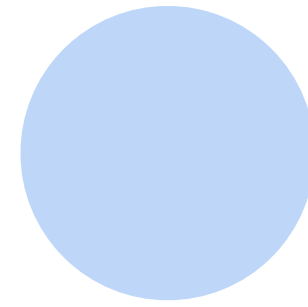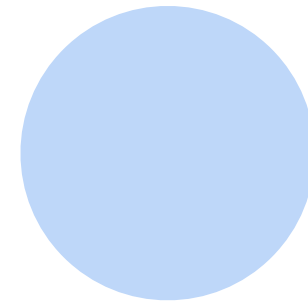Hidden layer
m neurons

# Multi-layer perceptron: Matrix Notation

- Input $\mathbf{x} \in \mathbb{R}^d$
- Hidden $\mathbf{W}^{(1)} \in \mathbb{R}^{m \times d}, \mathbf{b} \in \mathbb{R}^m$
- Intermediate output

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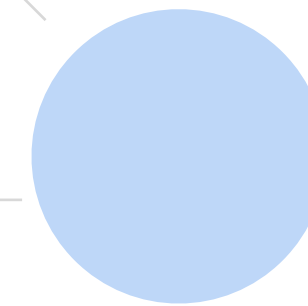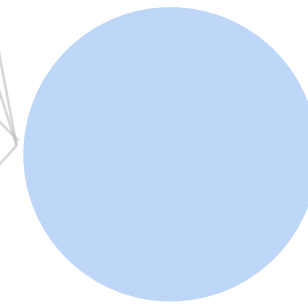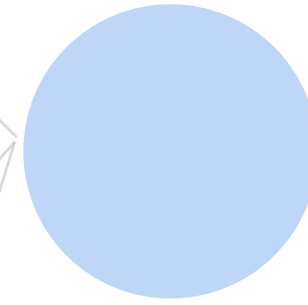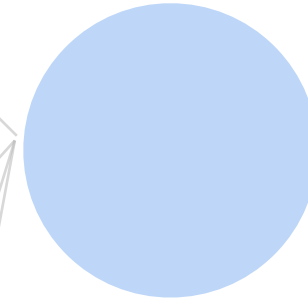
$$\mathbf{h} \in \mathbb{R}^m$$

Input

Hidden layer
m neurons

**Classify cats vs. dogs**

Input

Hidden layer
100 neurons

Output

# Multi-layer perceptron

Input

Hidden layer
m neurons

Output

Why do we need an a nonlinear activation?

# Multi-layer perceptron

Input

Hidden layer
m neurons

Output

Why do we need an a nonlinear activation?

$$\mathbf{h} = \mathbf{W}\mathbf{x} + \mathbf{b}$$

$$f = \mathbf{w}_2^T \mathbf{h} + b_2$$

hence $f = \mathbf{w}_2^\top \mathbf{W}\mathbf{x} + b'$

# Neural network for K-way classification

- K outputs in the final layer



Input

Hidden layer
m=3 neurons

$\mathbf{x} \in \mathbb{R}^d$

$x_1$

$x_2$

$h_1 = \sigma(\sum_{i=1}^{d} x_i w_{1i}^{(1)} + b_1)$

$h_2 = \sigma(\sum_{i=1}^{d} x_i w_{2i}^{(1)} + b_2)$

$h_3 = \sigma(\sum_{i=1}^{d} x_i w_{3i}^{(1)} + b_3)$

$w_{11}^{(2)}$

$w_{12}^{(2)}$

$w_{13}^{(2)}$

Output

$f_1 = \sum_{i=1}^{m} h_i w_{1i}^{(2)} + b_1'$

No activation function applied in output layer

# Neural network for K-way classification
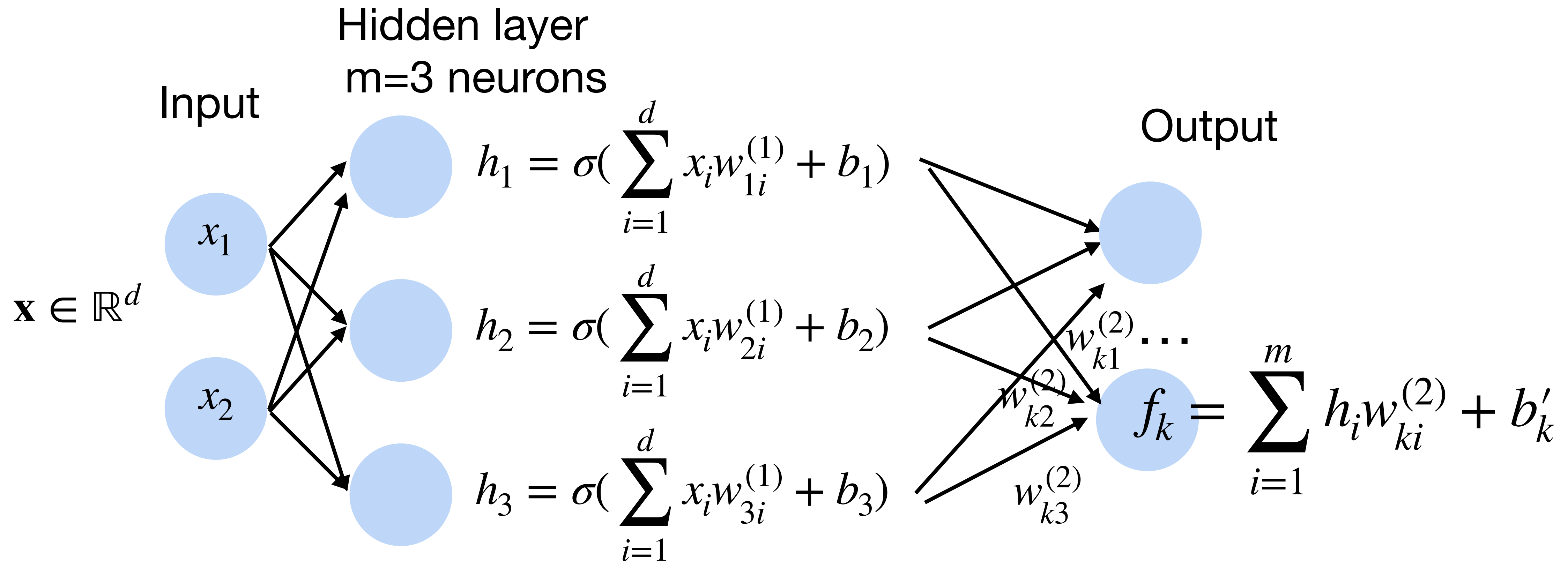
- K outputs units in the final layer

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

Input

Hidden layer
m=3 neurons

Output

$\mathbf{x} \in \mathbb{R}^d$

$x_1$

$x_2$

$h_1 = \sigma(\sum_{i=1}^{d} x_i w_{1i}^{(1)} + b_1)$

$h_2 = \sigma(\sum_{i=1}^{d} x_i w_{2i}^{(1)} + b_2)$

$h_3 = \sigma(\sum_{i=1}^{d} x_i w_{3i}^{(1)} + b_3)$

$w_{k1}^{(2)} \ldots$

$w_{k2}^{(2)}$

$w_{k3}^{(2)}$

$f_k = \sum_{i=1}^{m} h_i w_{ki}^{(2)} + b_k'$

# Softmax

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



Input

Hidden layer
m=3 neurons

Output

$\mathbf{x} \in \mathbb{R}^d$

$x_1$

$x_2$

$\cdots$

$$p(y \,|\, \mathbf{x}) = \text{softmax}(f)$$

$$= \frac{\exp f_y(x)}{\sum_{k=1}^{K} \exp f_k(x)}$$
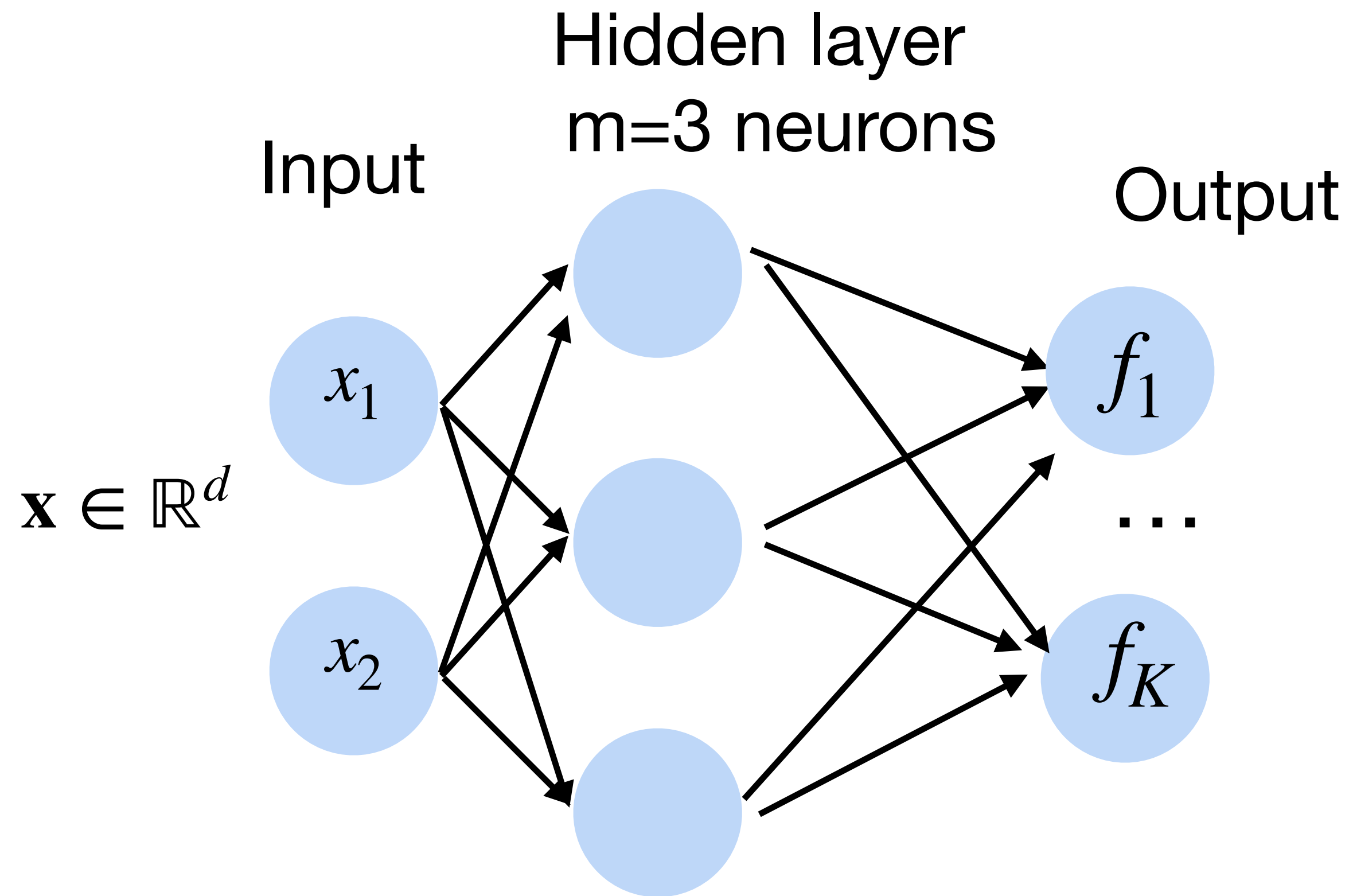
# **Softmax**

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



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

$$= \frac{\exp f_y(x)}{\sum_{k=1}^{K} \exp f_k(x)}$$

# Softmax

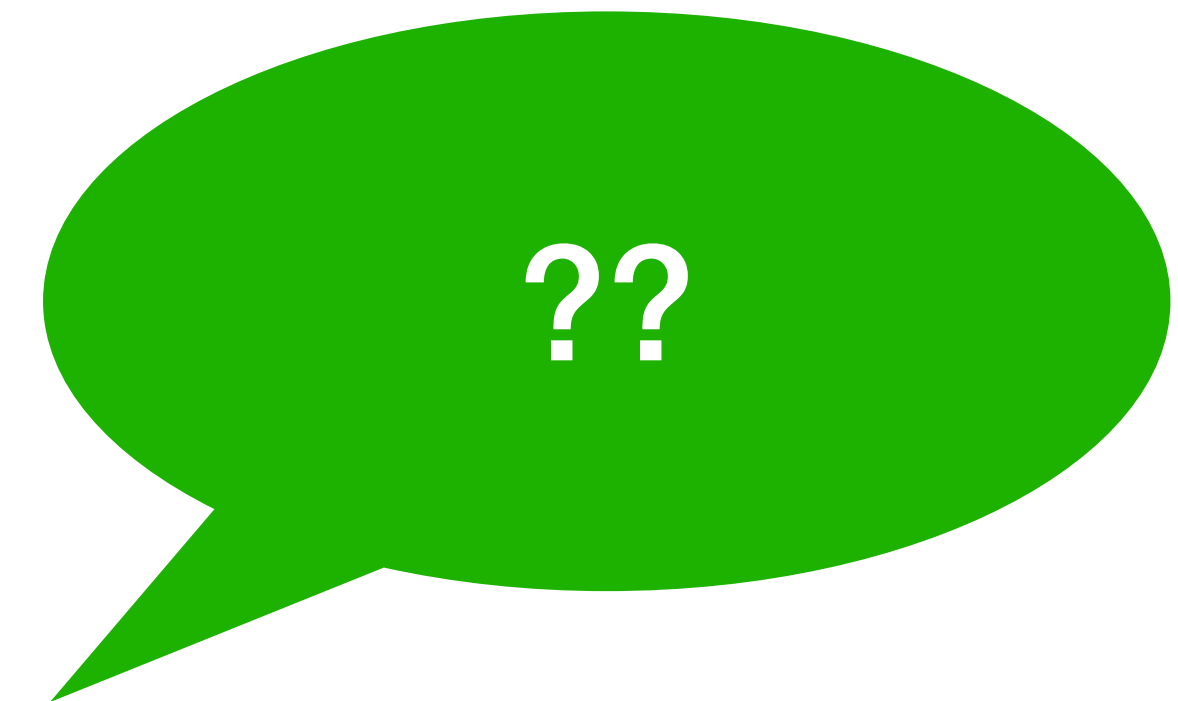Turns outputs f into probabilities (sum up to 1 across K classes)

# 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 brid
```
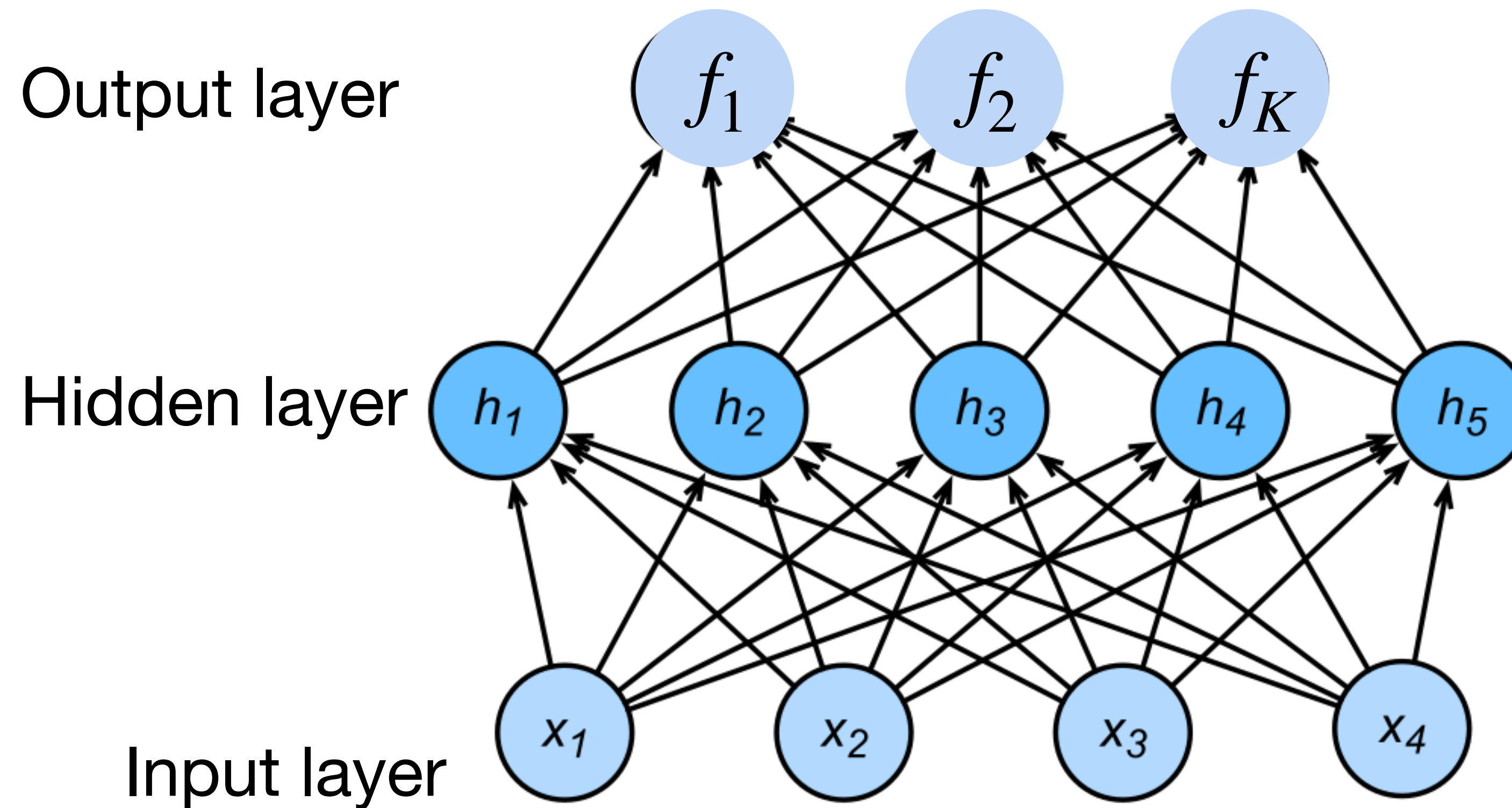
https://www.kaggle.com/c/human-protein-atlas-image-classification

# More complicated neural networks

Output layer $f_1$ $f_2$ $f_K$

Hidden layer $h_1$ $h_2$ $h_3$ $h_4$ $h_5$

Input layer $x_1$ $x_2$ $x_3$ $x_4$

# More complicated neural networks

$$p_1, p_2, \ldots, p_K = \text{softmax}(f_1, f_2, \ldots, f_K)$$



Output layer

$f_1$  $f_2$  $f_K$

Hidden layer

$h_1$  $h_2$  $h_3$  $h_4$  $h_5$

Input layer

$x_1$  $x_2$  $x_3$  $x_4$

# More complicated neural networks

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

$$\mathbf{h} = \sigma(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b})$$
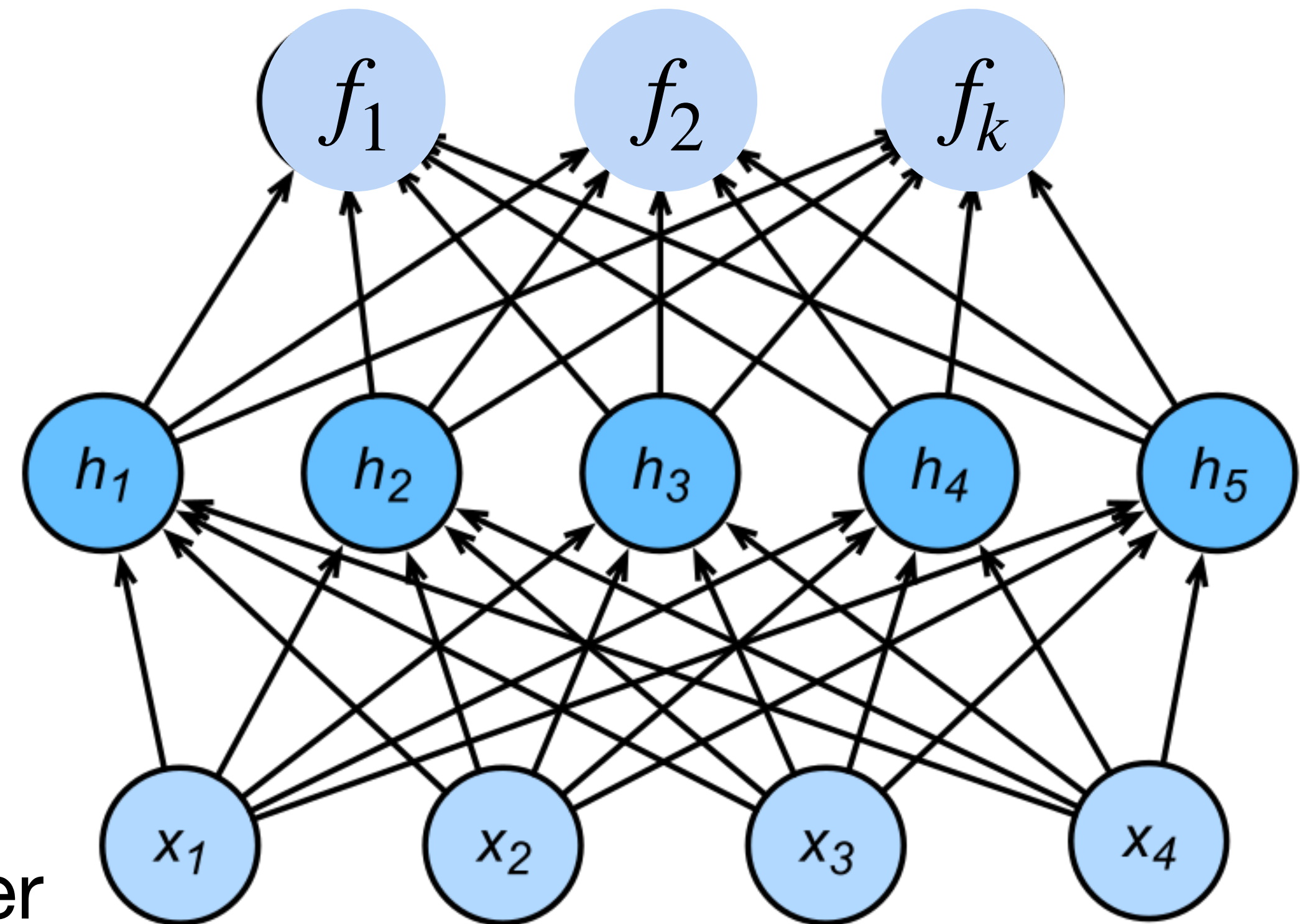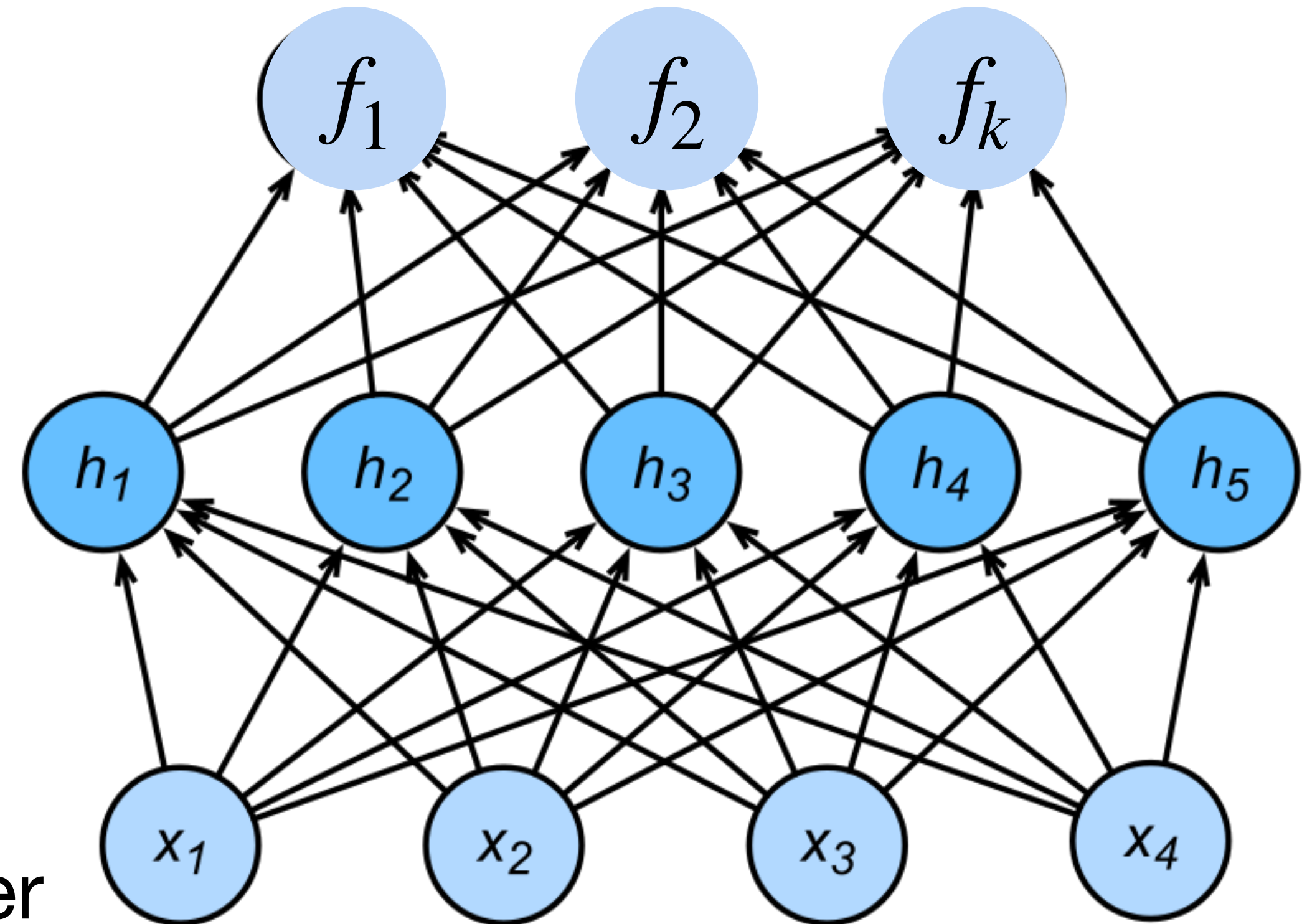
$$\mathbf{f} = \sigma(\mathbf{W}^{(2)}\mathbf{h} + \mathbf{b}^{(2)})$$

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

Output layer

Hidden layer

Input layer

# More complicated neural networks

- Input $\mathbf{x} \in \mathbb{R}^d$

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

$$p_1, p_2, \ldots, p_K = \text{softmax}(f_1, f_2, \ldots, f_K)$$

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

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

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

Output layer

$f_1$  $f_2$  $f_k$

Hidden layer

$h_1$  $h_2$  $h_3$  $h_4$  $h_5$

Input layer

$x_1$  $x_2$  $x_3$  $x_4$

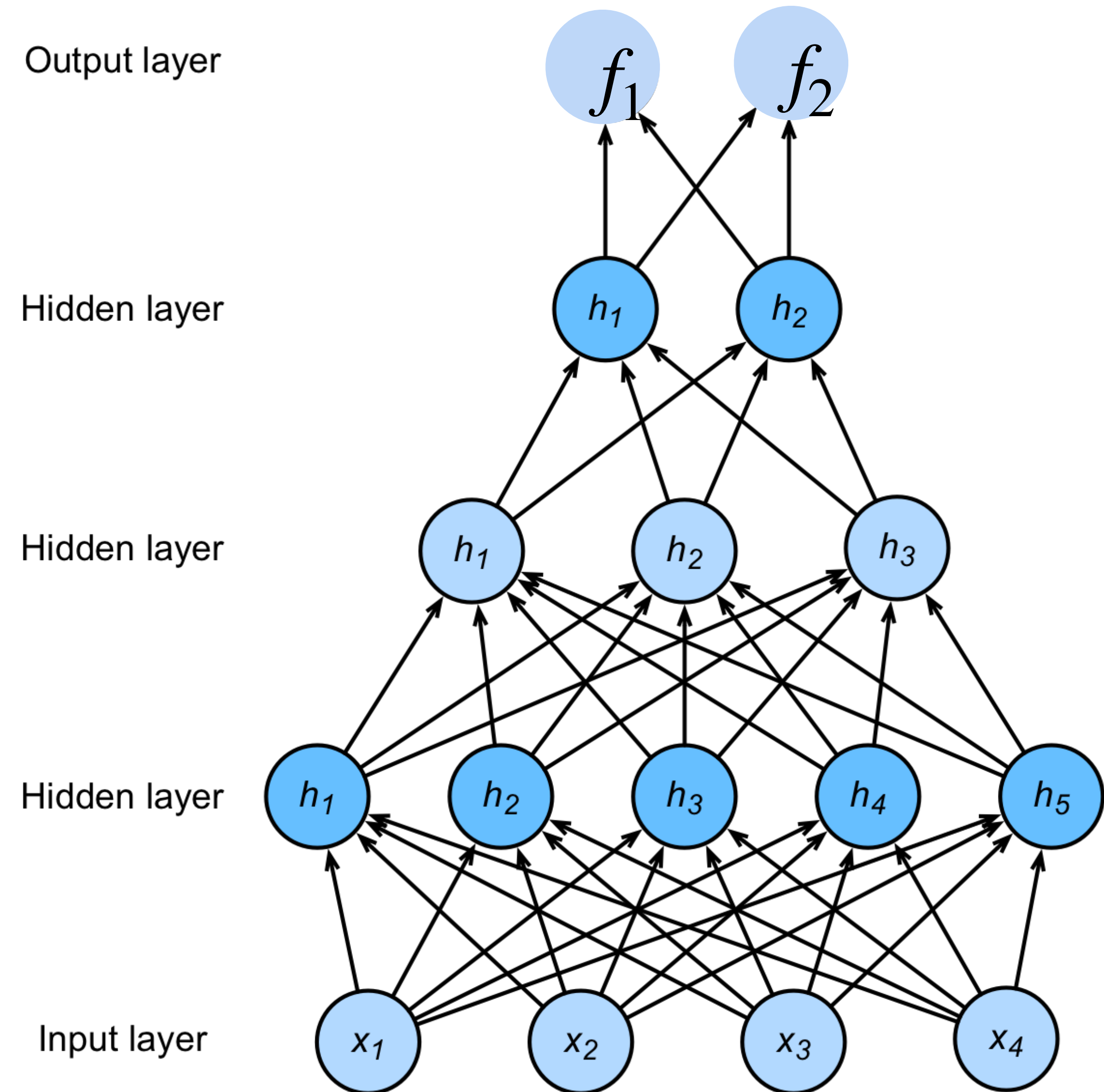# More complicated neural networks: multiple hidden layers

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

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

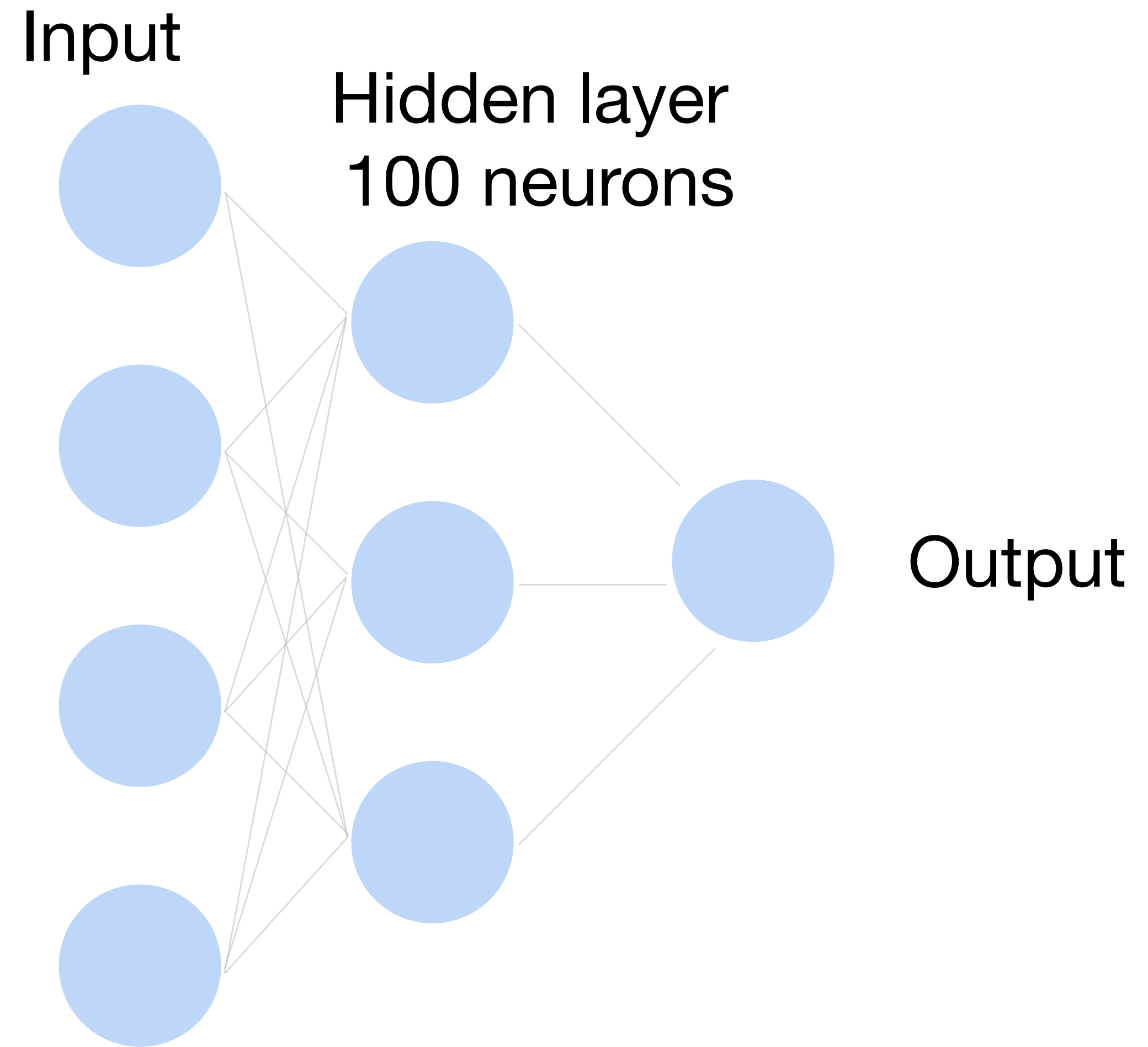$$\mathbf{h}_3 = \sigma(\mathbf{W}_3\mathbf{h}_2 + \mathbf{b}_3)$$

$$\mathbf{f} = \mathbf{W}_4\mathbf{h}_3 + \mathbf{b}_4$$

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

# How to train a neural network?

**Classify cats vs. dogs**

Input

Hidden layer
100 neurons

Output
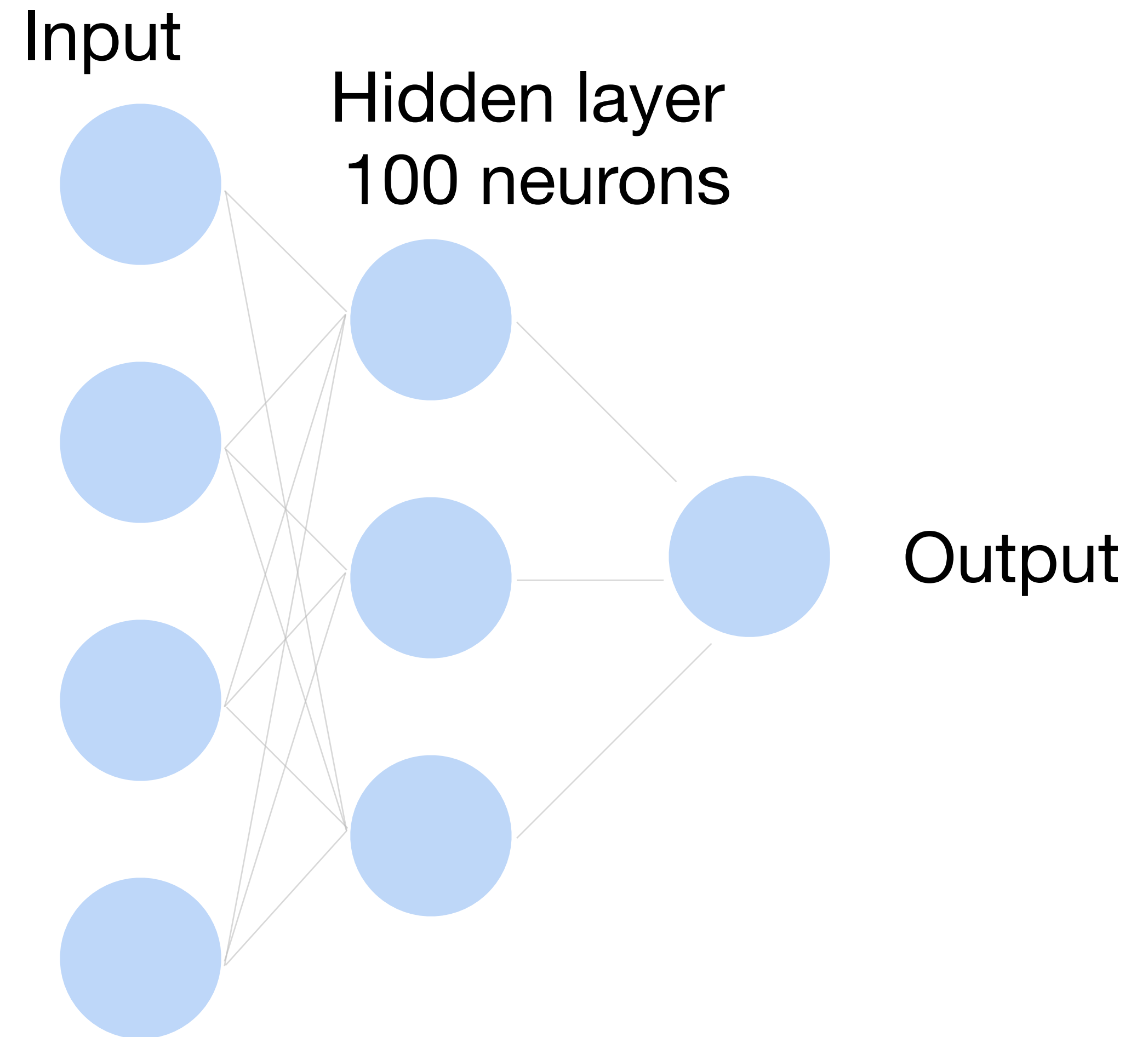
# How to train a neural network? Binary classification

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

$\hat{y}$ Model output for example $\mathbf{x}$
(This is a function of all weights W)
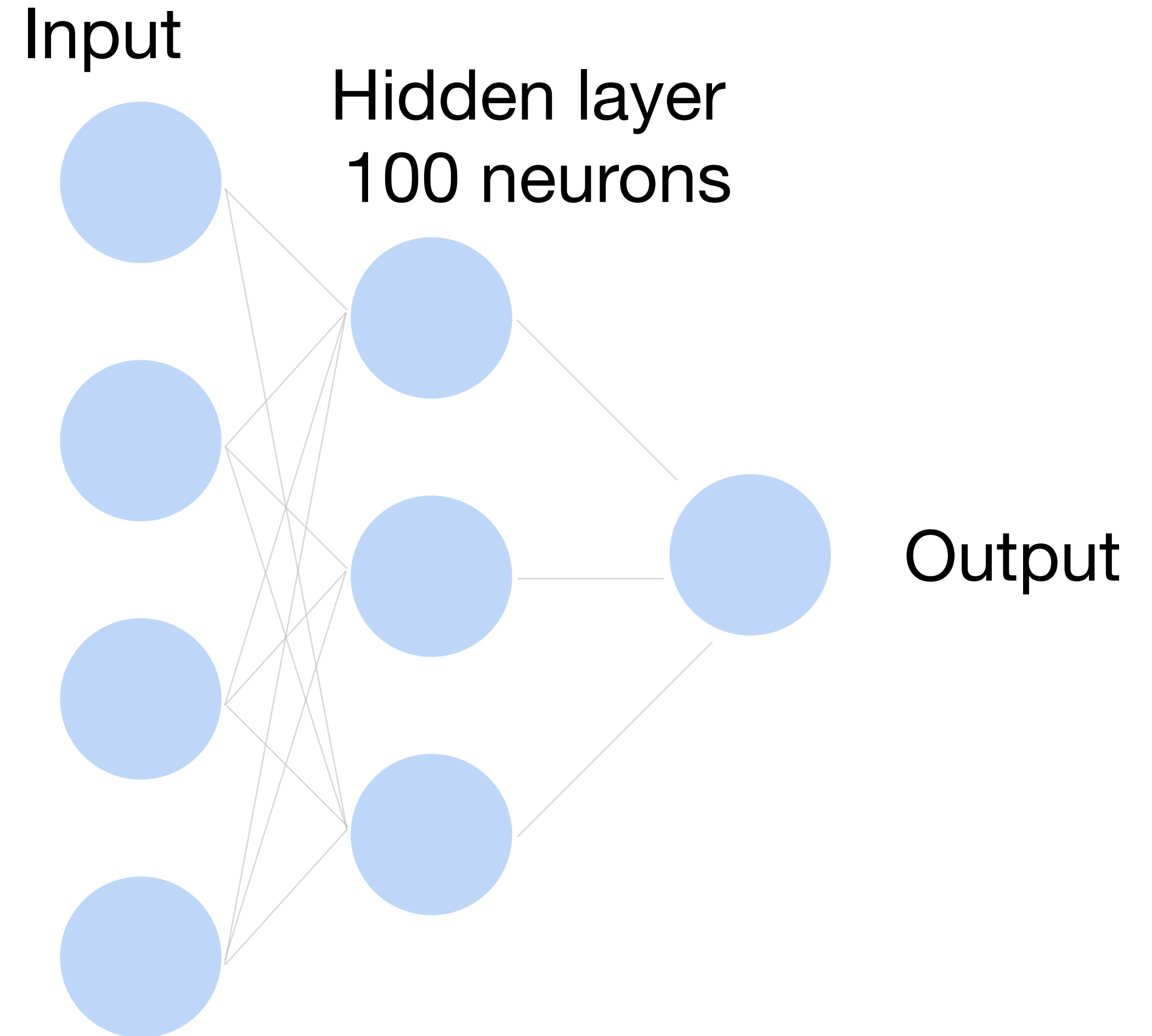
$y$ Ground truth label for example $\mathbf{x}$

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

**We want** $\hat{y} \to 1$ **when** $y = 1$,
**and** $\hat{y} \to 0$ **when** $y = 0$

Input

Hidden layer
100 neurons

Output

# How to train a neural network? Binary classification

**Loss function:** $\dfrac{1}{|D|} \displaystyle\sum_{(\mathbf{x},y)\in D} \ell(\mathbf{x},y)$



Input

Hidden layer
100 neurons

Output

# How to train a neural network? Binary classification

**Loss function:** $\dfrac{1}{|D|} \displaystyle\sum_{(\mathbf{x},y)\in D} \ell(\mathbf{x}, y)$

**Per-sample loss:**

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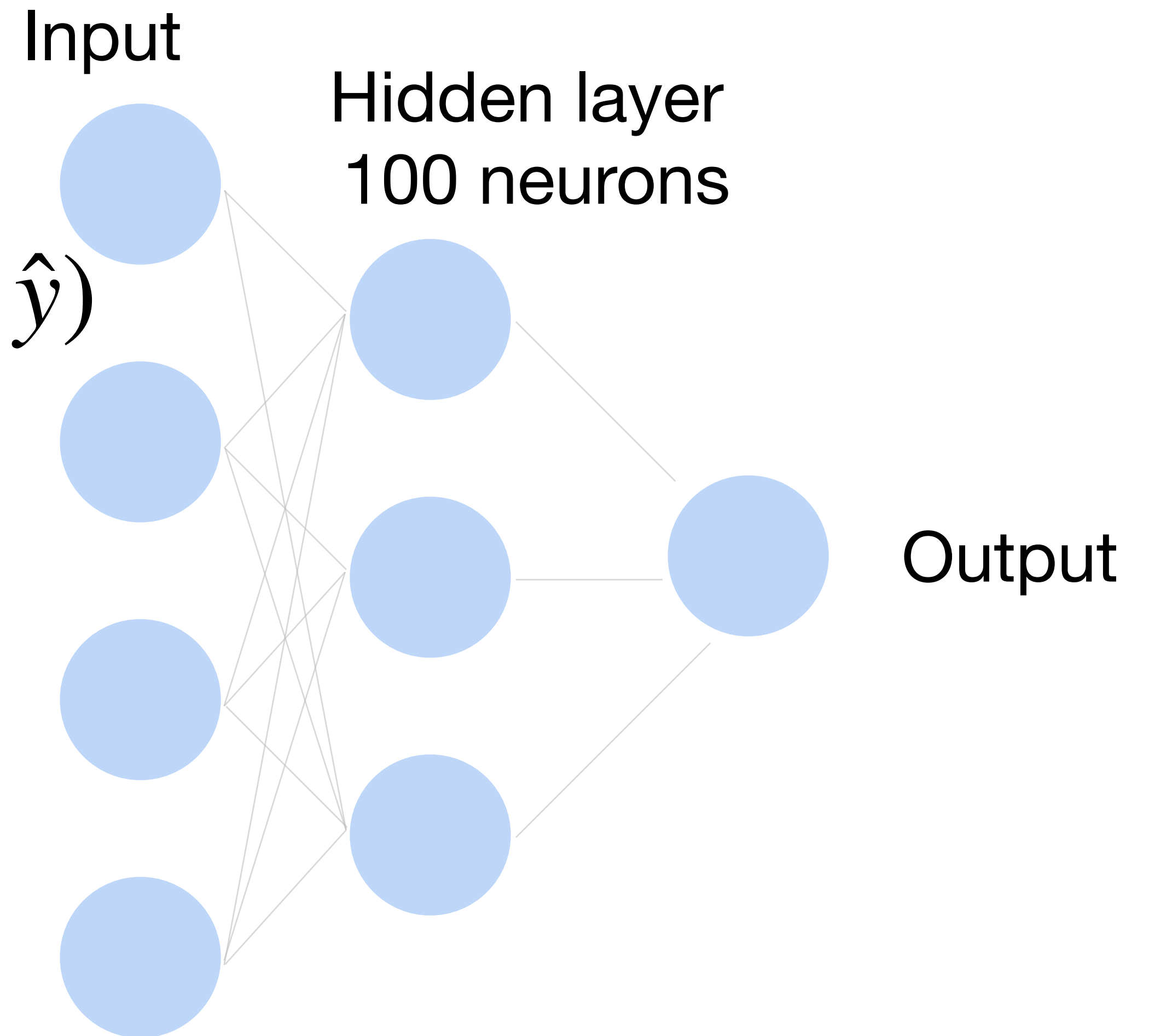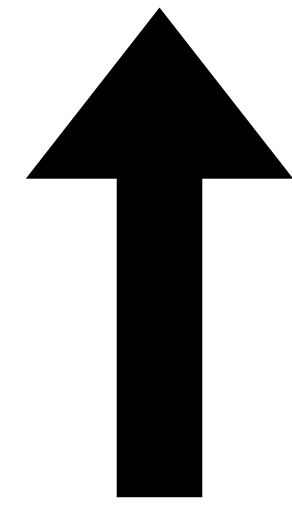
Input

Hidden layer
100 neurons

Output

# How to train a neural network? Binary classification

**Loss function:** $\dfrac{1}{|D|} \displaystyle\sum_{(\mathbf{x},y) \in D} \ell(\mathbf{x}, y)$
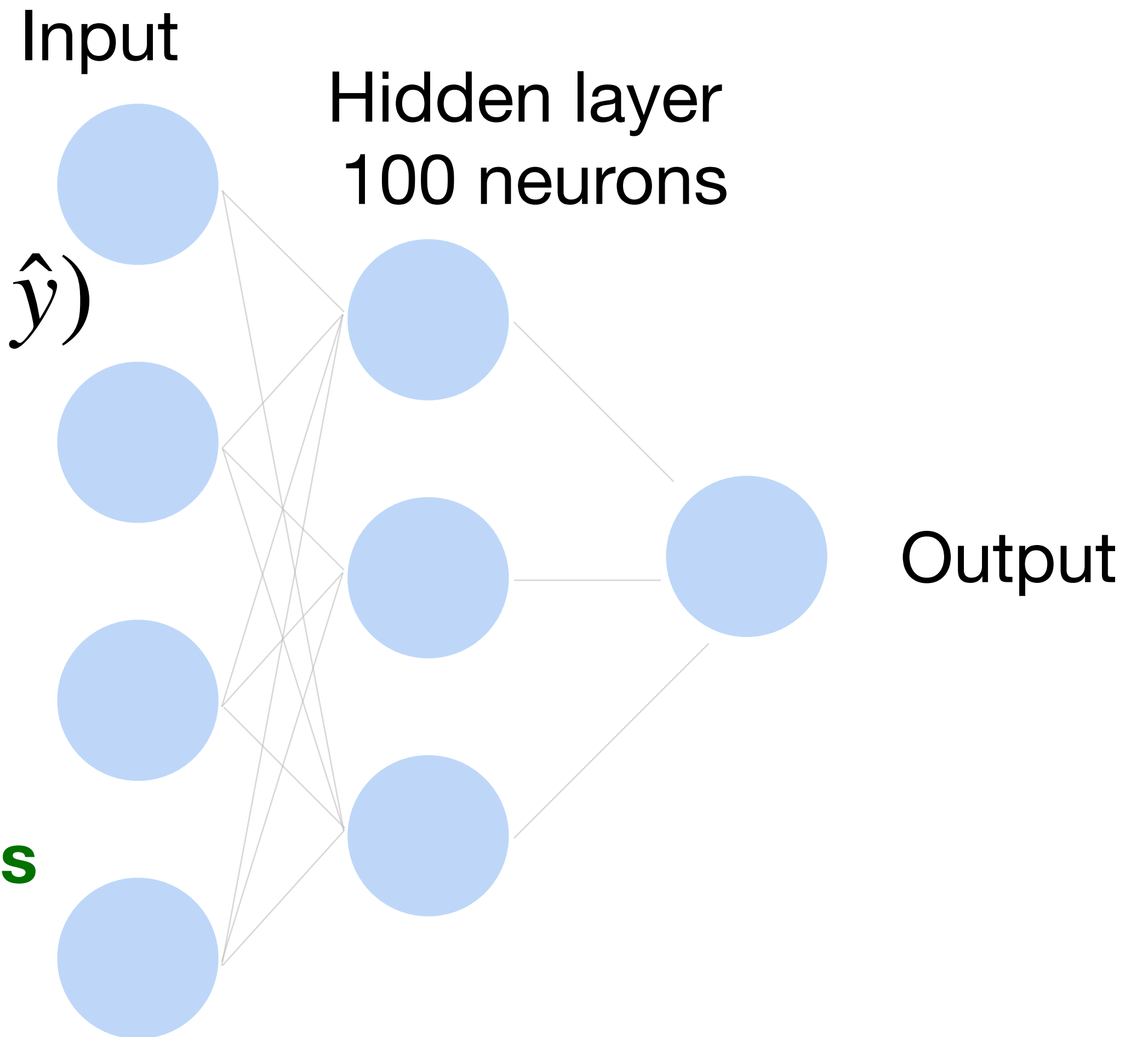
Input

Hidden layer
100 neurons

**Per-sample loss:**

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

Output

**Negative log likelihood**
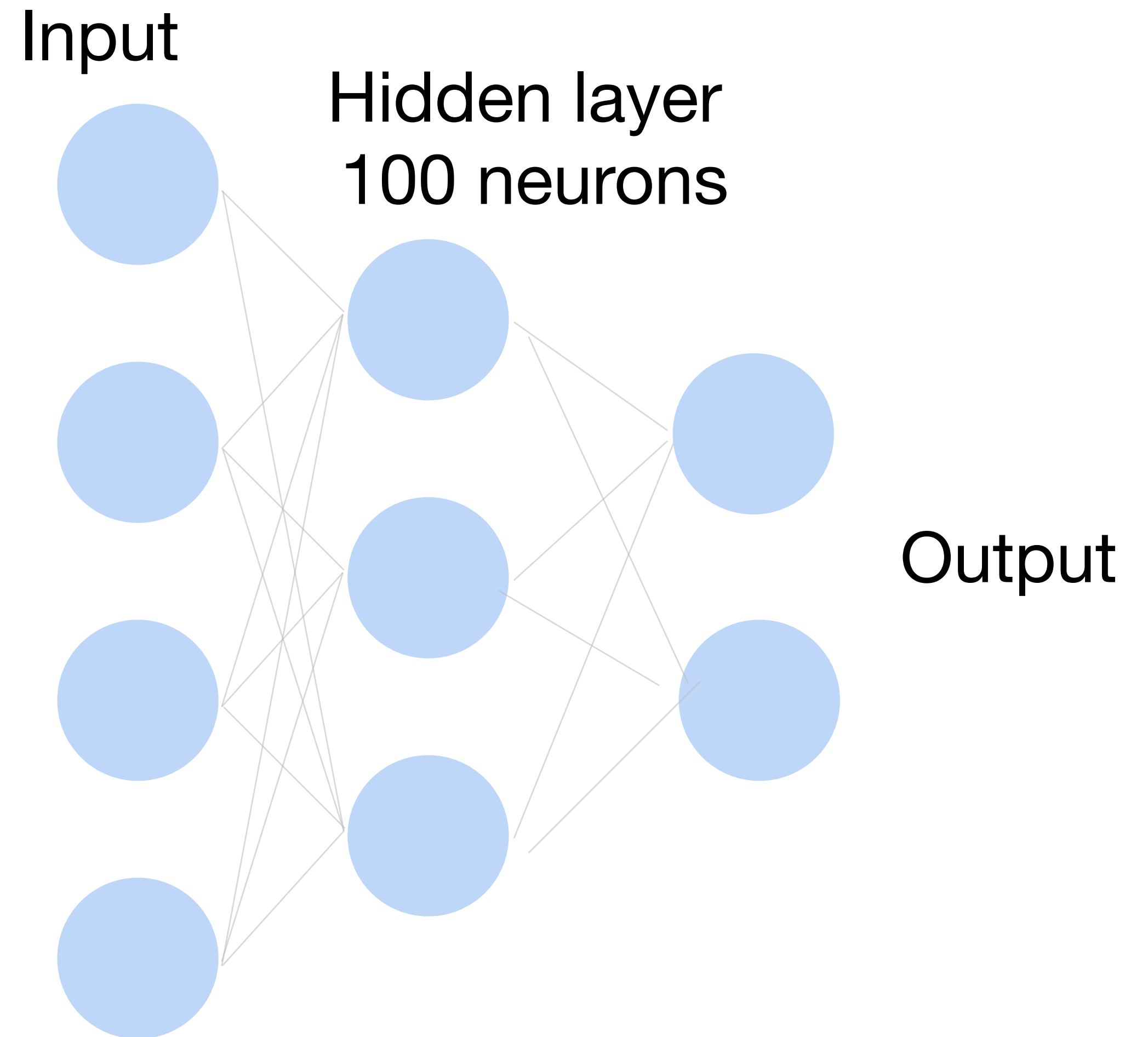
**Also known as binary cross-entropy loss**

# How to train a neural network? Multiclass

**Loss function:** $\dfrac{1}{|D|} \displaystyle\sum_{(\mathbf{x},y)\in D} \ell(\mathbf{x}, y)$

Input

Hidden layer
100 neurons

Output

$Y$

$y$

# How to train a neural network? Multiclass

**Loss function:** $\dfrac{1}{|D|} \displaystyle\sum_{(\mathbf{x},y)\in D} \ell(\mathbf{x}, y)$

**Per-sample loss:**

$$\ell(\mathbf{x}, y) = \sum_{k=1}^{K} -Y_k \log p_k = -\log p_y$$

where $Y$ is one-hot encoding of $y$

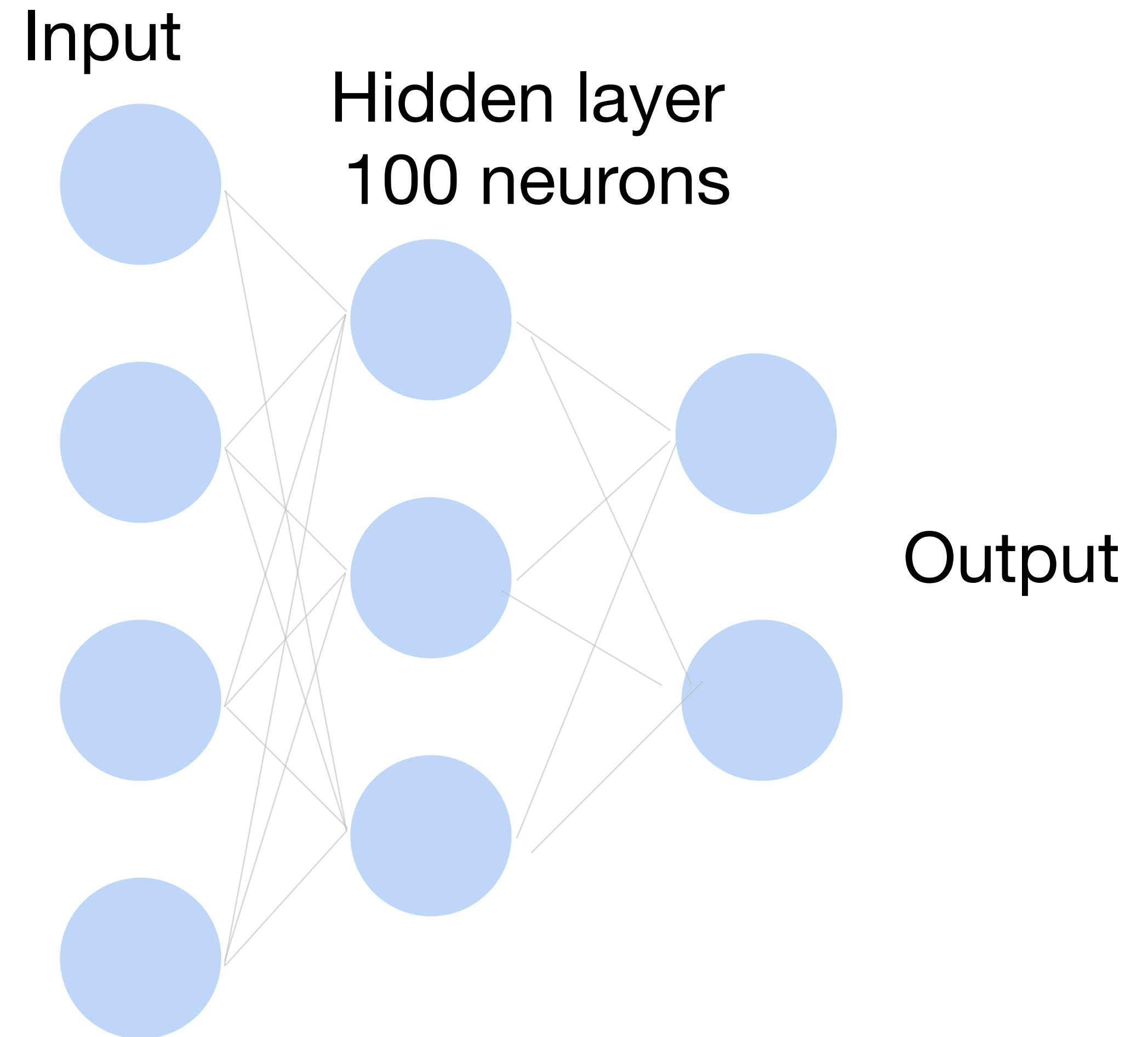Input

Hidden layer
100 neurons
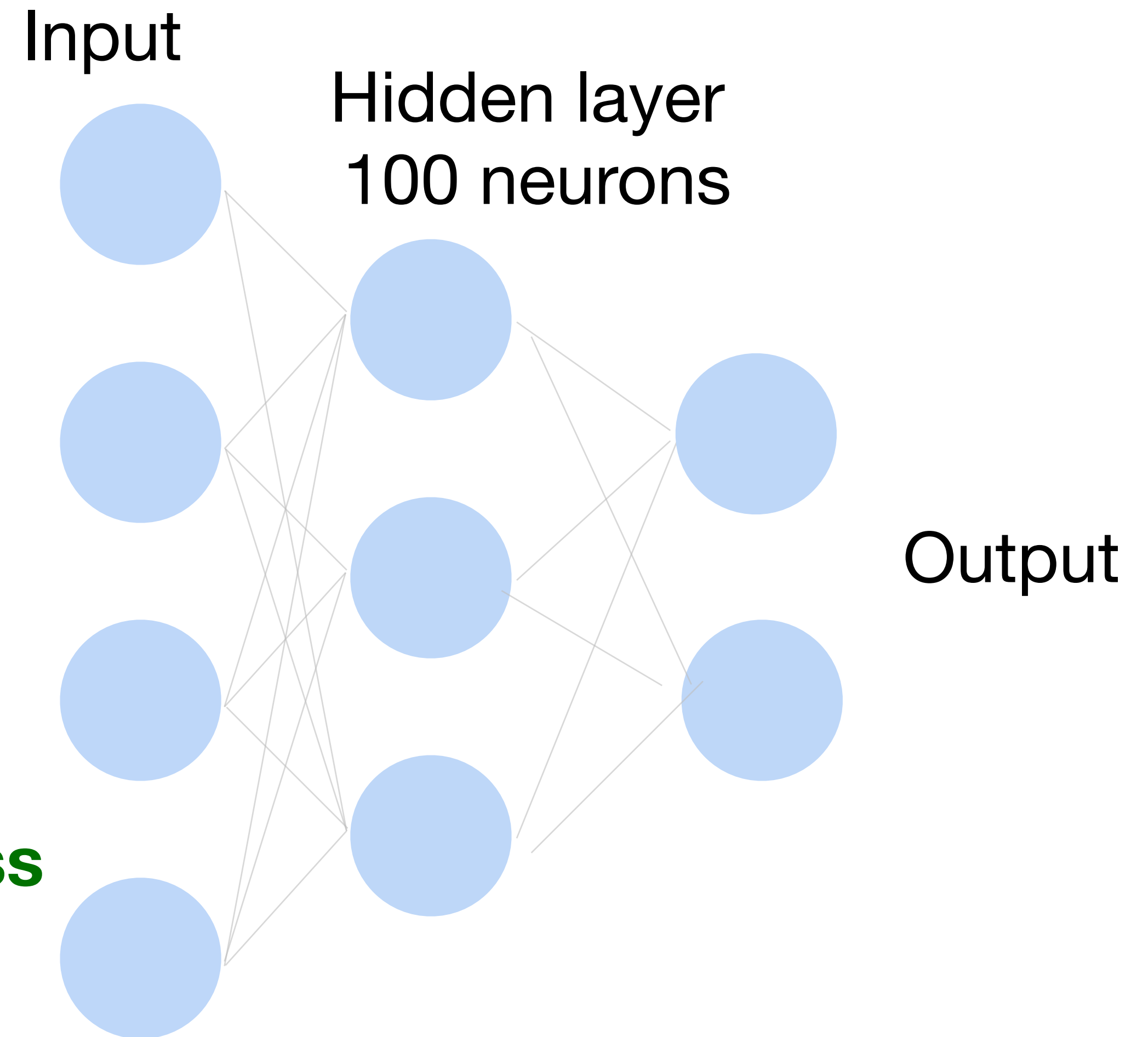
Output

# How to train a neural network? Multiclass

**Loss function:** $\dfrac{1}{|D|} \displaystyle\sum_{(\mathbf{x}, y) \in D} \ell(\mathbf{x}, y)$

**Per-sample loss:**

$$\ell(\mathbf{x}, y) = \sum_{k=1}^{K} - Y_k \log p_k = -\log p_y$$

where $Y$ is one-hot encoding of $y$

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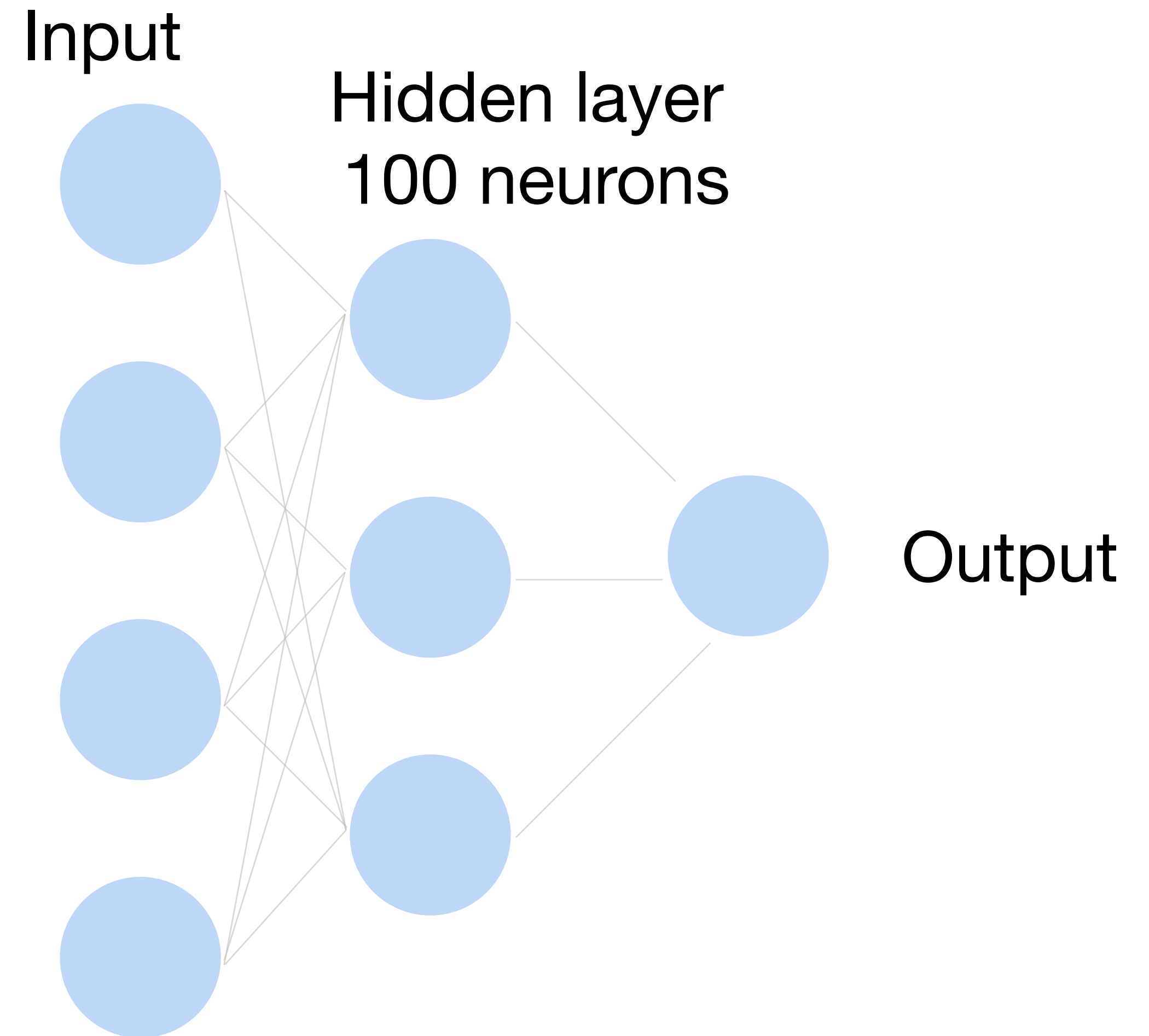
Input

Hidden layer
100 neurons

Output

# How to train a neural network?

Update the weights W to minimize the loss function

$$L = \frac{1}{|D|} \sum_{(\mathbf{x},y) \in D} \ell(\mathbf{x}, y)$$

**Use gradient descent!**
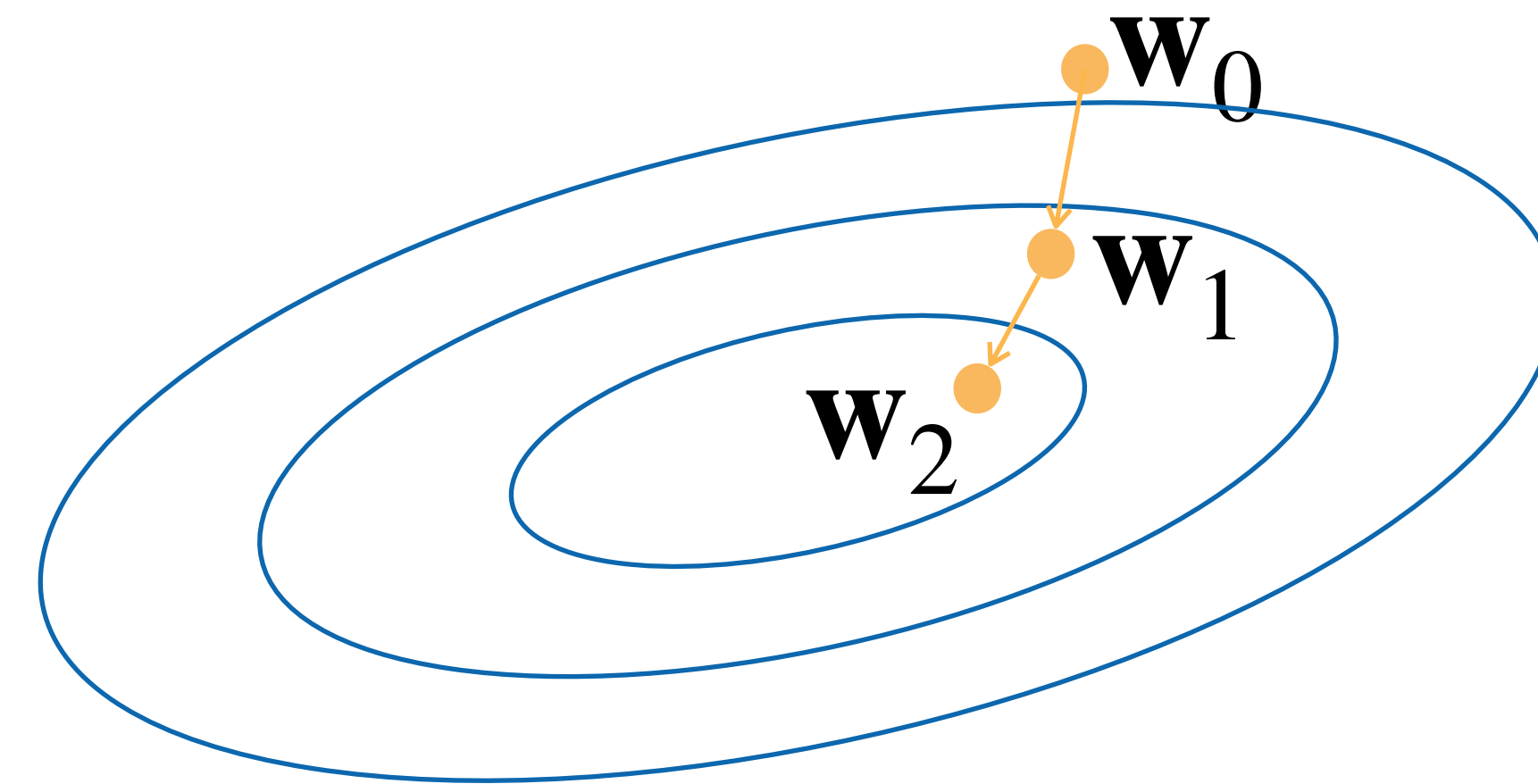
Input

Hidden layer
100 neurons

Output

# Gradient Descent



- Choose a learning rate $\alpha > 0$
- Initialize the model parameters $w_0$
- For t =1,2,…
    - Update parameters:

$$\mathbf{w}_t = \mathbf{w}_{t-1} - \alpha\frac{\partial L}{\partial \mathbf{w}_{t-1}}$$

$$= \mathbf{w}_{t-1} - \alpha\frac{1}{|D|}\sum_{(\mathbf{x},y)\in D}\frac{\partial \ell(\mathbf{x}, y)}{\partial \mathbf{w}_{t-1}}$$

- Repeat until converges

# Gradient Descent



- Choose a learning rate $\alpha > 0$
- Initialize the model parameters $w_0$
- For t =1,2,…

    - Update parameters:

$$\mathbf{w}_t = \mathbf{w}_{t-1} - \alpha \frac{\partial L}{\partial \mathbf{w}_{t-1}}$$

$$= \mathbf{w}_{t-1} - \alpha \frac{1}{|D|} \sum_{(\mathbf{x}, y) \in D} \frac{\partial \ell(\mathbf{x}, y)}{\partial \mathbf{w}_{t-1}}$$
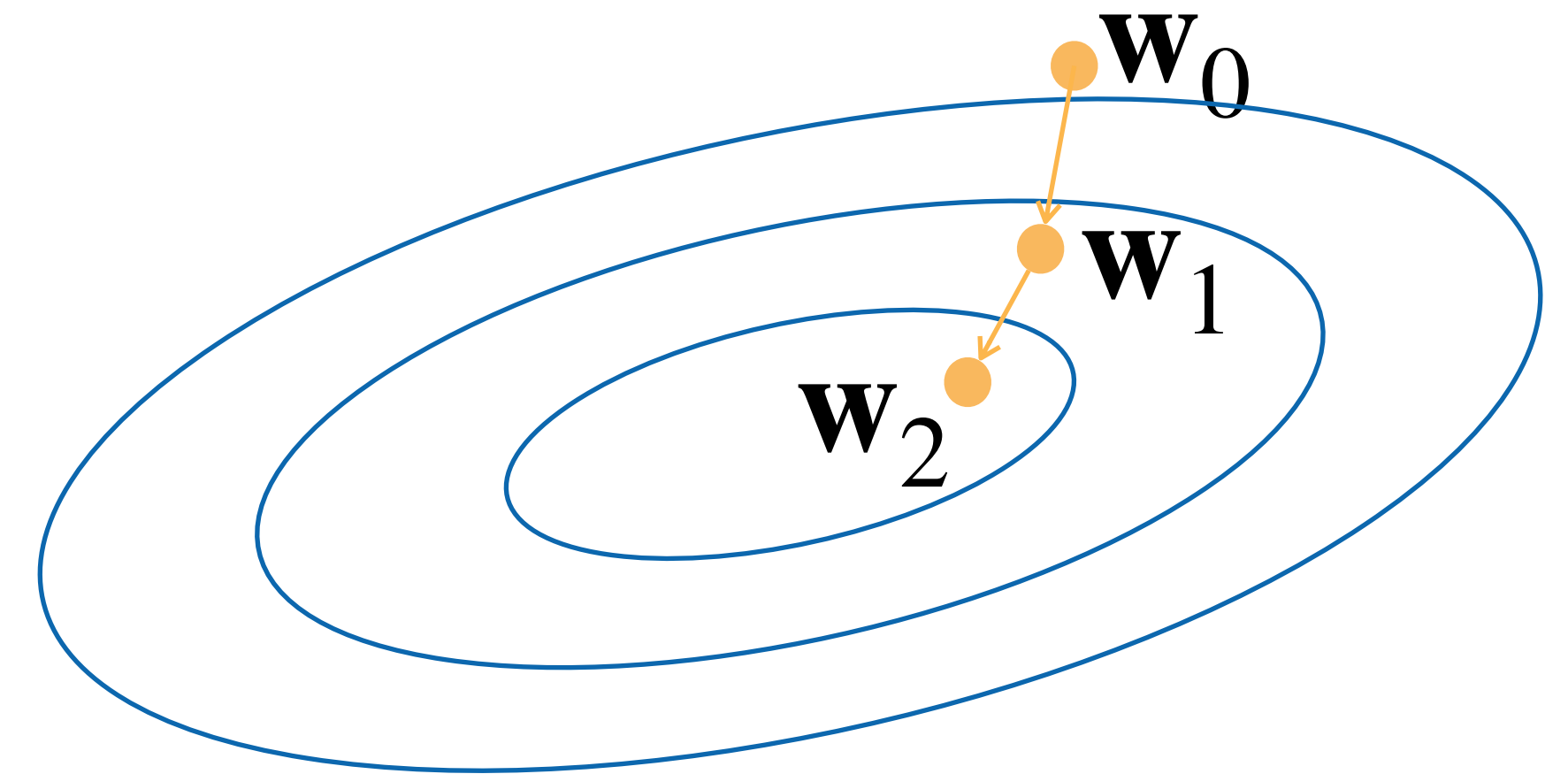
D can be very large. Expensive

- Repeat until converges

# Gradient Descent



- Choose a learning rate $\alpha > 0$
- Initialize the model parameters $w_0$
- For t =1,2,…

  - Update parameters:

$$\mathbf{w}_t = \mathbf{w}_{t-1} - \alpha \frac{\partial L}{\partial \mathbf{w}_{t-1}}$$

$$= \mathbf{w}_{t-1} - \alpha \frac{1}{|D|} \sum_{(\mathbf{x},y) \in D} \frac{\partial \ell(\mathbf{x}, y)}{\partial \mathbf{w}_{t-}}$$
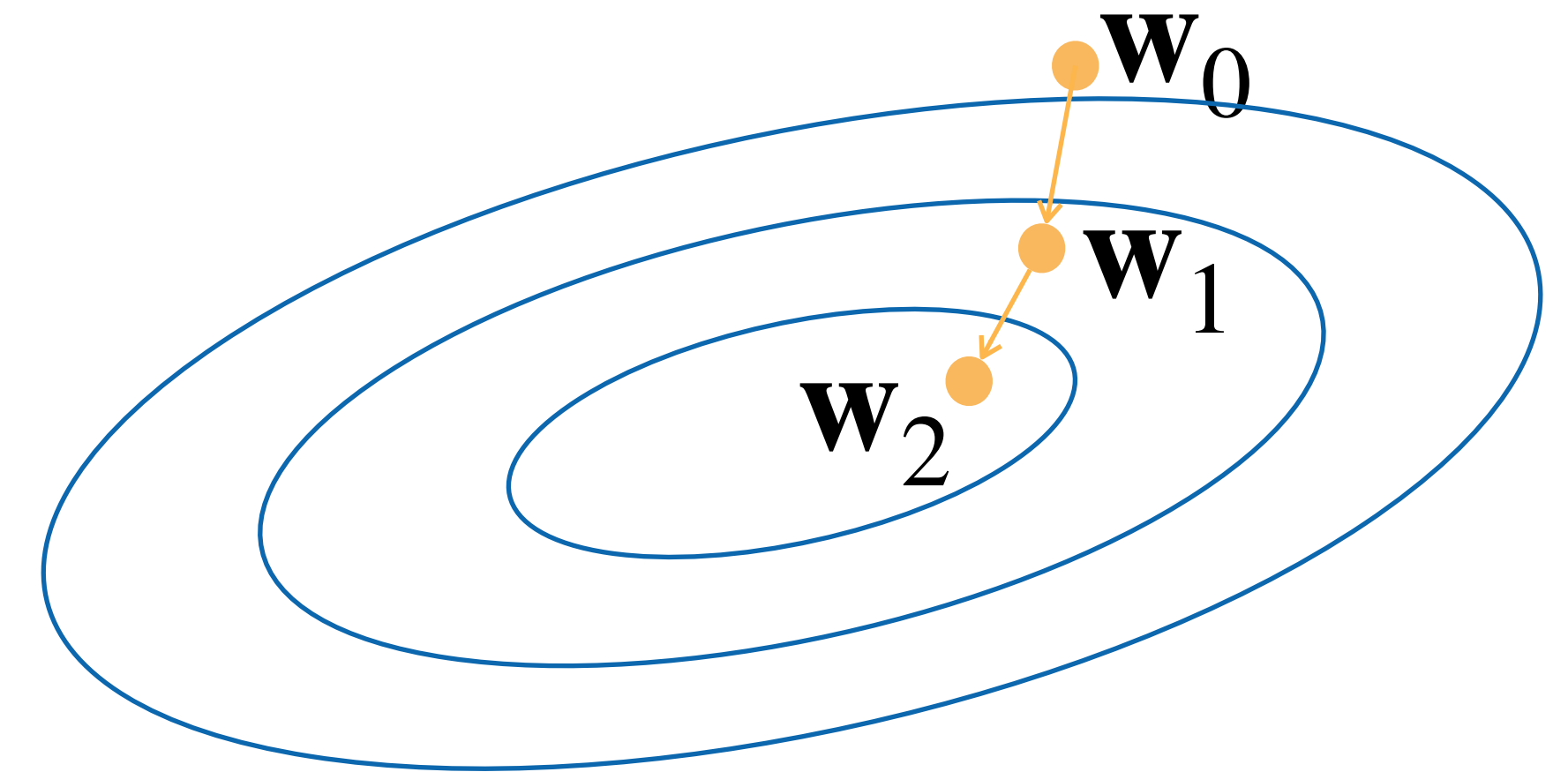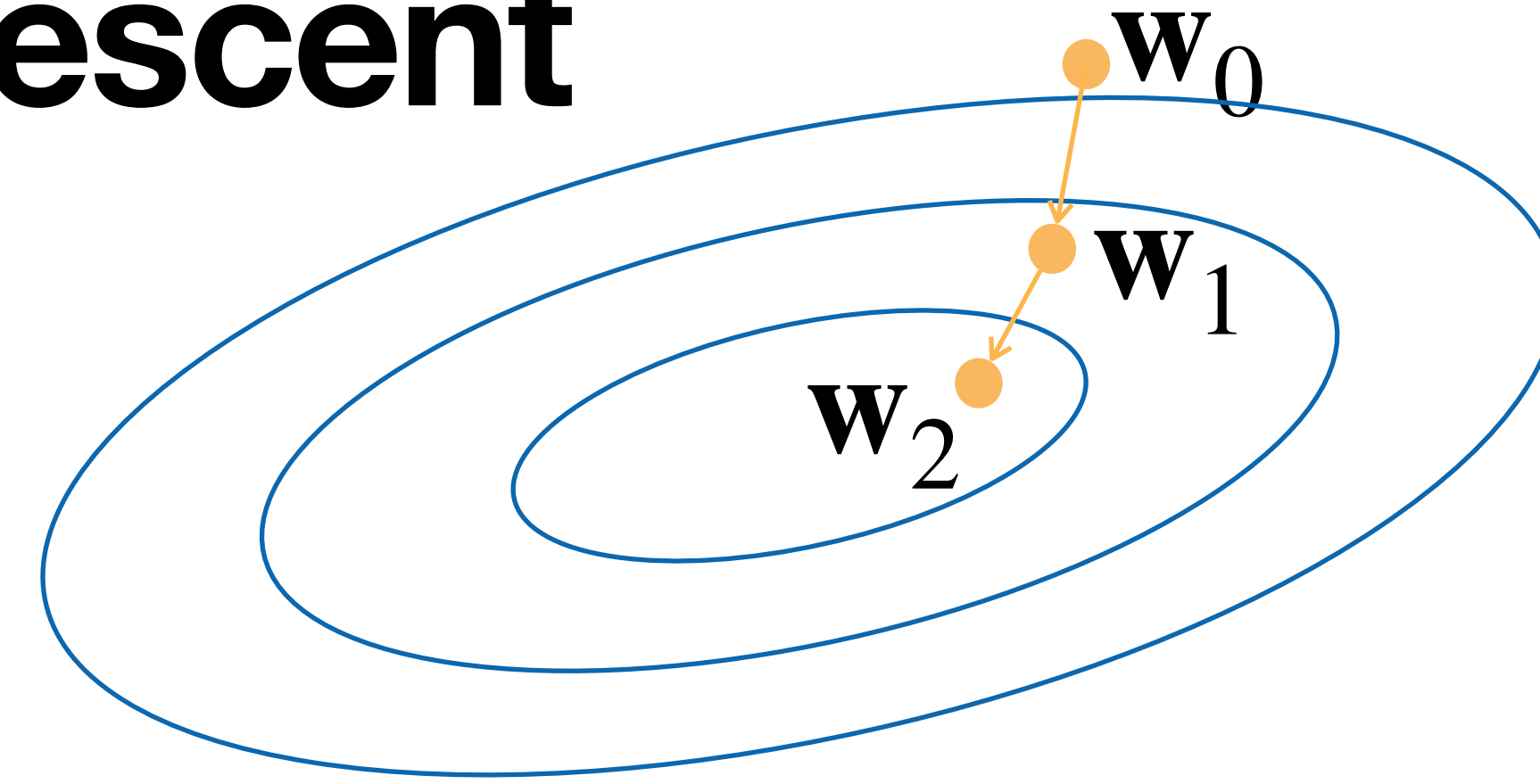
D can be very large. Expensive

The gradient w.r.t. all parameters is obtained by concatenating the partial derivatives w.r.t. each parameter

- Repeat until converges

# Minibatch Stochastic Gradient Descent



- Choose a learning rate $\alpha > 0$
- Initialize the model parameters $w_0$
- For t =1,2,…

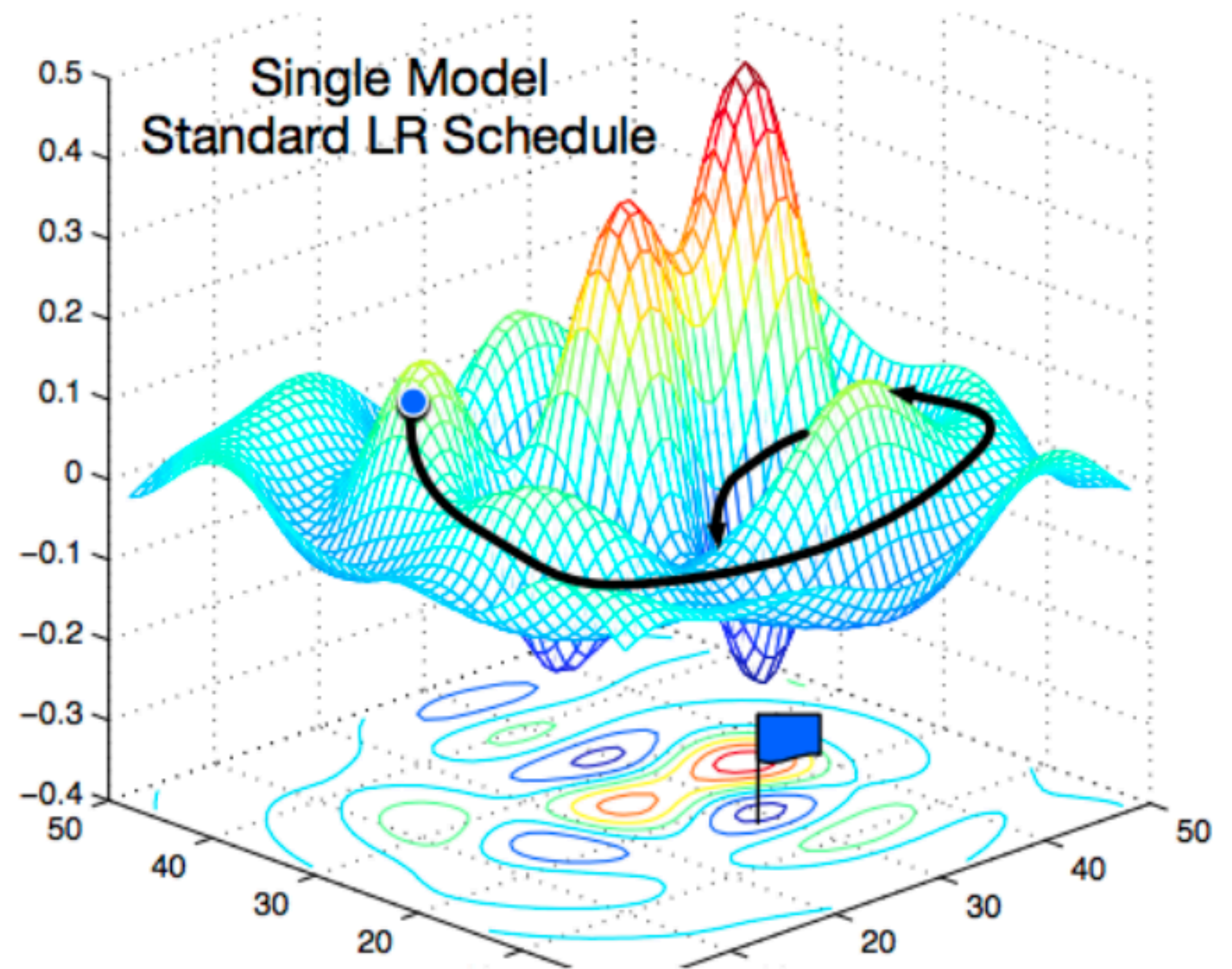  - **Randomly sample a subset (mini-batch)** $B \subset D$
    Update parameters:

$$\mathbf{w}_t = \mathbf{w}_{t-1} - \alpha \frac{1}{|B|} \sum_{(\mathbf{x},y) \in B} \frac{\partial \ell(\mathbf{x}, y)}{\partial \mathbf{w}_{t-1}}$$

- Repeat until converges

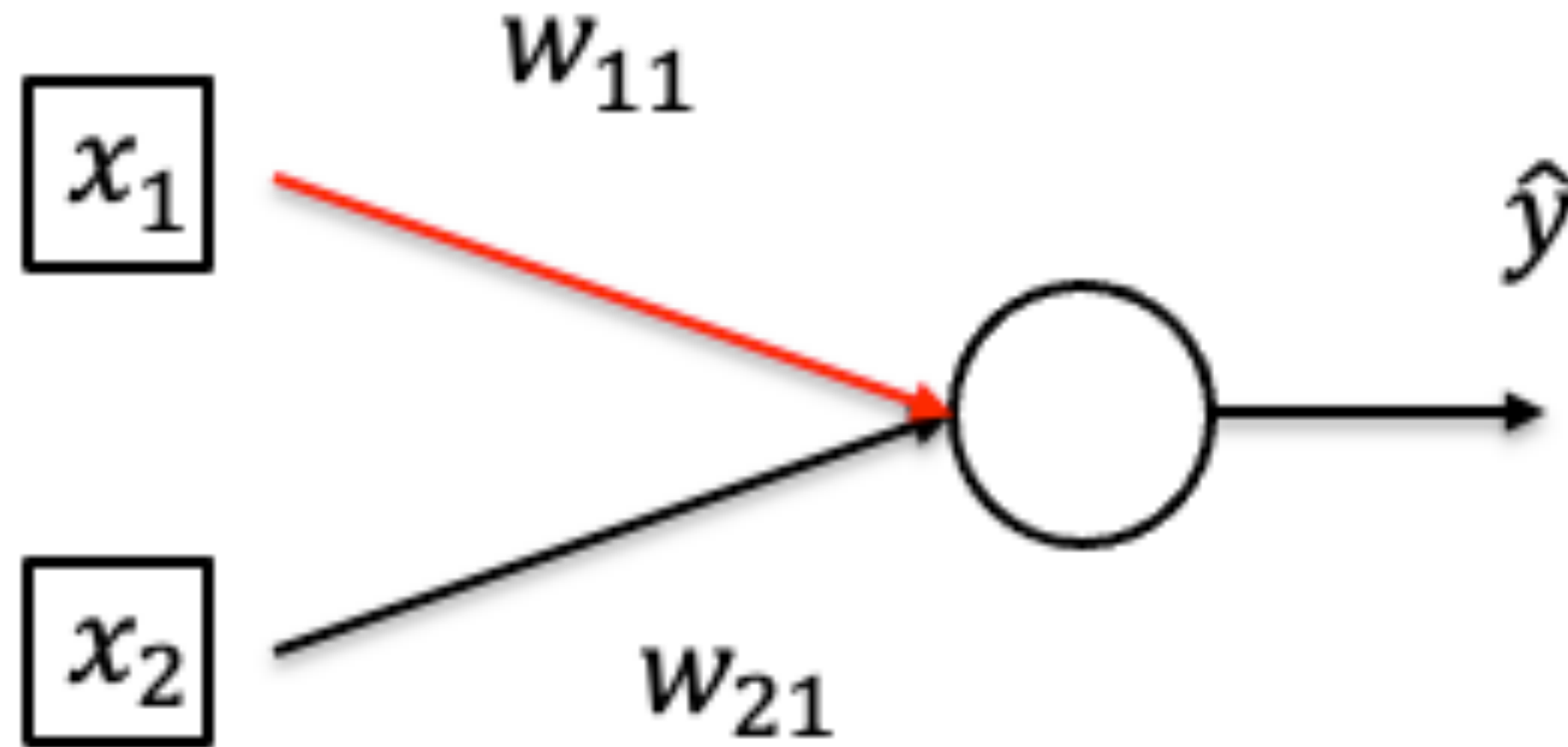# Non-convex Optimization
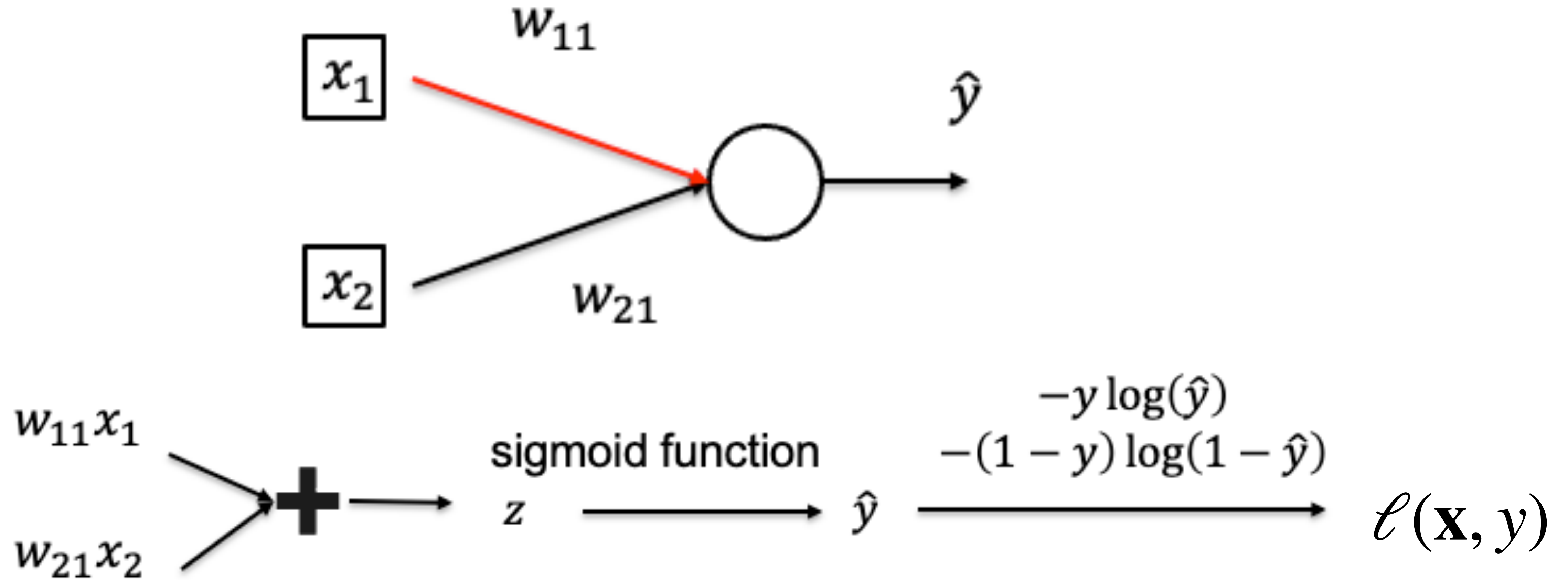


Single Model
Standard LR Schedule

[Gao and Li et al., 2018]

# Calculate Gradient (on one data point)



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

# Calculate Gradient (on one data point)

# Calculate Gradient (on one data point)



$$w_{11}x_1$$
$$w_{21}x_2$$

$$\text{sigmoid function}$$

$$z \longrightarrow \hat{y}$$

$$\frac{\partial \hat{y}}{\partial z} = \sigma'(z)$$

$$-y \log(\hat{y})$$
$$-(1-y)\log(1-\hat{y})$$

$$\ell(\mathbf{x}, y)$$

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

# Calculate Gradient (on one data point)



$$w_{11}x_1$$
$$w_{21}x_2$$

$$z \xrightarrow{\text{sigmoid function}} \hat{y} \xrightarrow{\substack{-y\log(\hat{y}) \\ -(1-y)\log(1-\hat{y})}} \ell(\mathbf{x}, y)$$

$$\frac{\partial \hat{y}}{\partial z} = \sigma'(z)$$

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

- By chain rule:

$$\frac{\partial l}{\partial w_{11}} = \frac{\partial l}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial w_{11}}$$
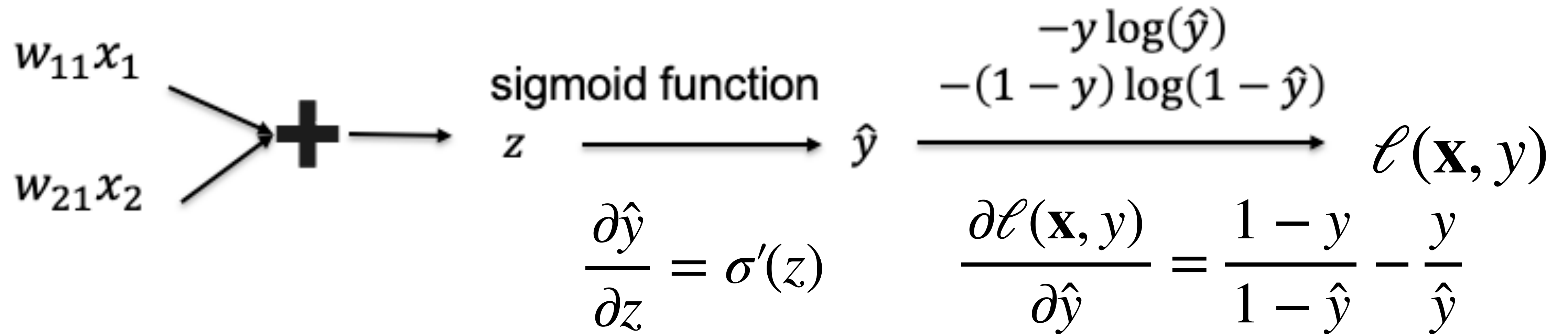
# Calculate Gradient (on one data point)

$w_{11}$

$x_1$

$w_{21}$

$x_2$

$\hat{y}$

$w_{11}x_1$

$w_{21}x_2$

**+**

sigmoid function

$z \longrightarrow \hat{y}$

$-y\log(\hat{y})$
$-(1-y)\log(1-\hat{y})$

$\longrightarrow \ell(\mathbf{x}, y)$

$$\frac{\partial \hat{y}}{\partial z} = \sigma'(z)$$

$$\frac{\partial \ell(\mathbf{x}, y)}{\partial \hat{y}} = \frac{1-y}{1-\hat{y}} - \frac{y}{\hat{y}}$$
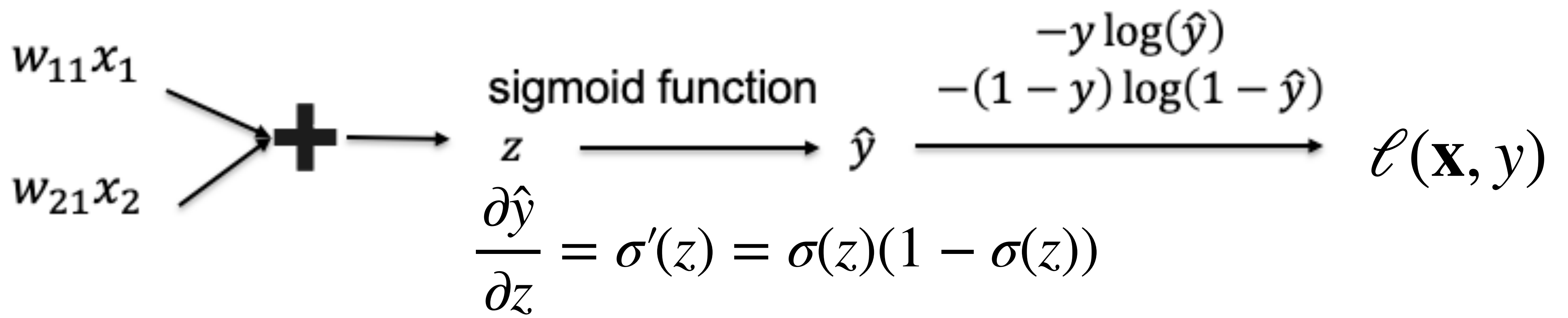
- By chain rule:

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

# Calculate Gradient (on one data point)
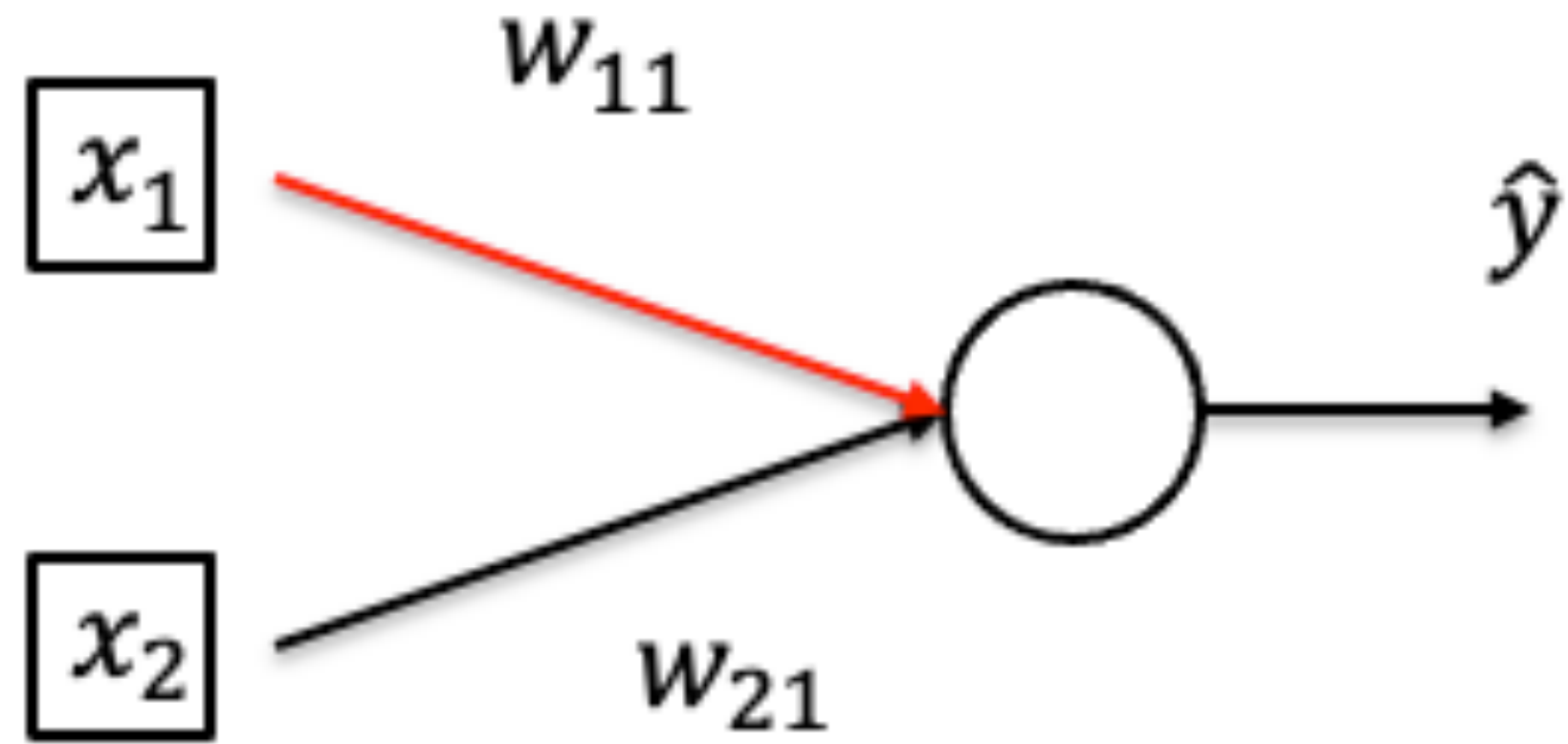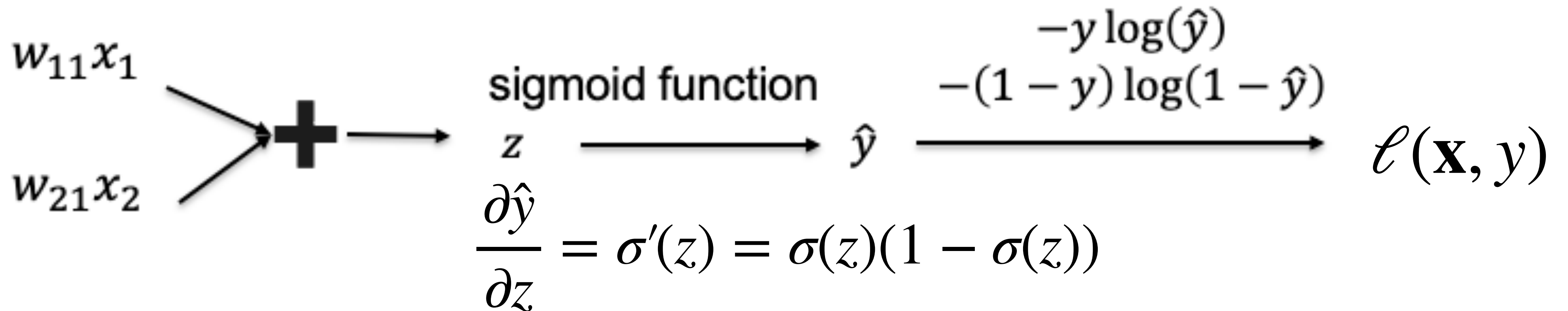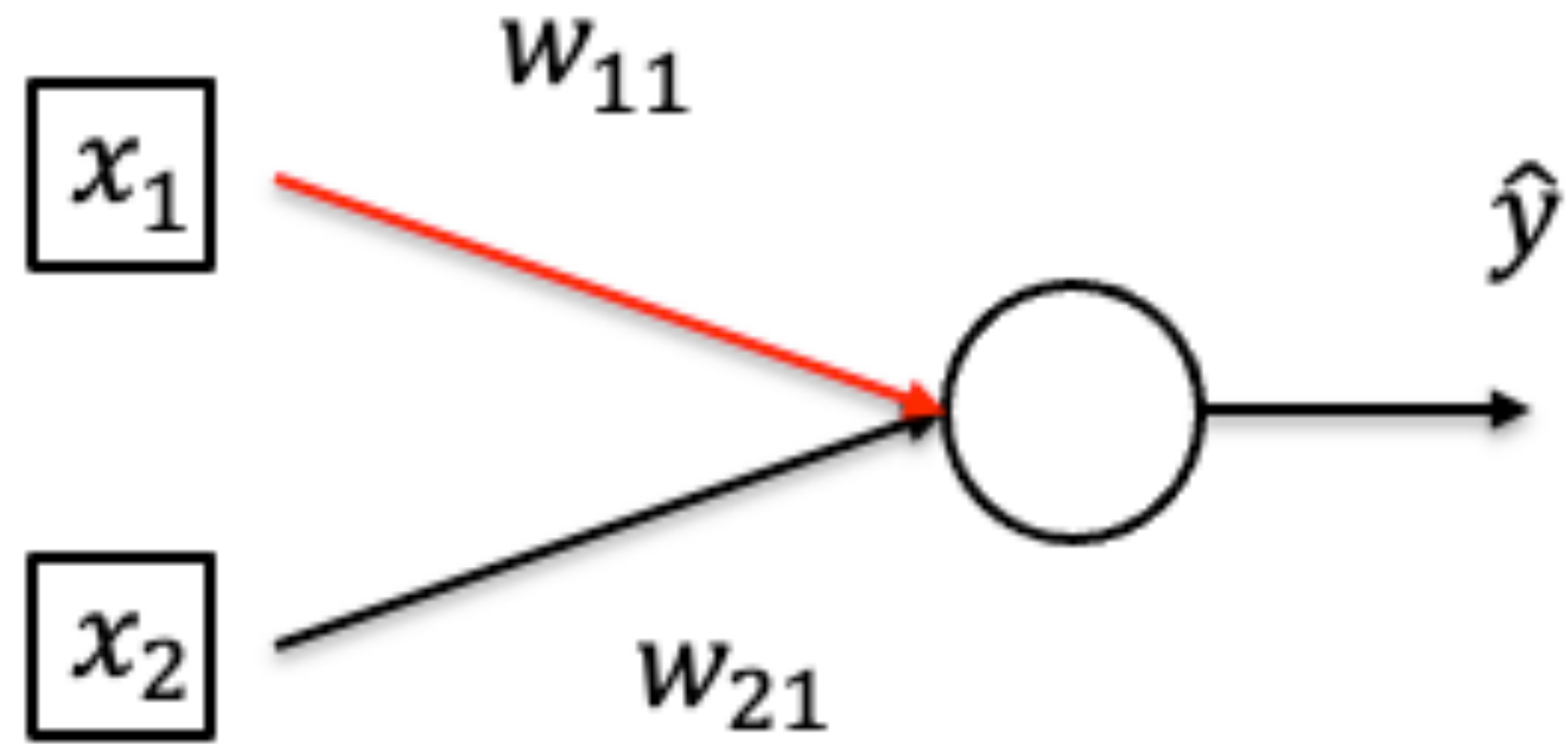


$$\frac{\partial \hat{y}}{\partial z} = \sigma'(z) = \sigma(z)(1 - \sigma(z))$$

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

# Calculate Gradient (on one data point)



$w_{11}x_1$

$w_{21}x_2$

$+$

sigmoid function

$z \longrightarrow \hat{y}$

$-y\log(\hat{y})$
$-(1-y)\log(1-\hat{y})$

$\ell(\mathbf{x}, y)$

$$\frac{\partial \hat{y}}{\partial z} = \sigma'(z) = \sigma(z)(1-\sigma(z))$$

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

# Calculate Gradient (on one data point)



$$\frac{\partial \hat{y}}{\partial z} = \sigma'(z) = \sigma(z)(1 - \sigma(z))$$
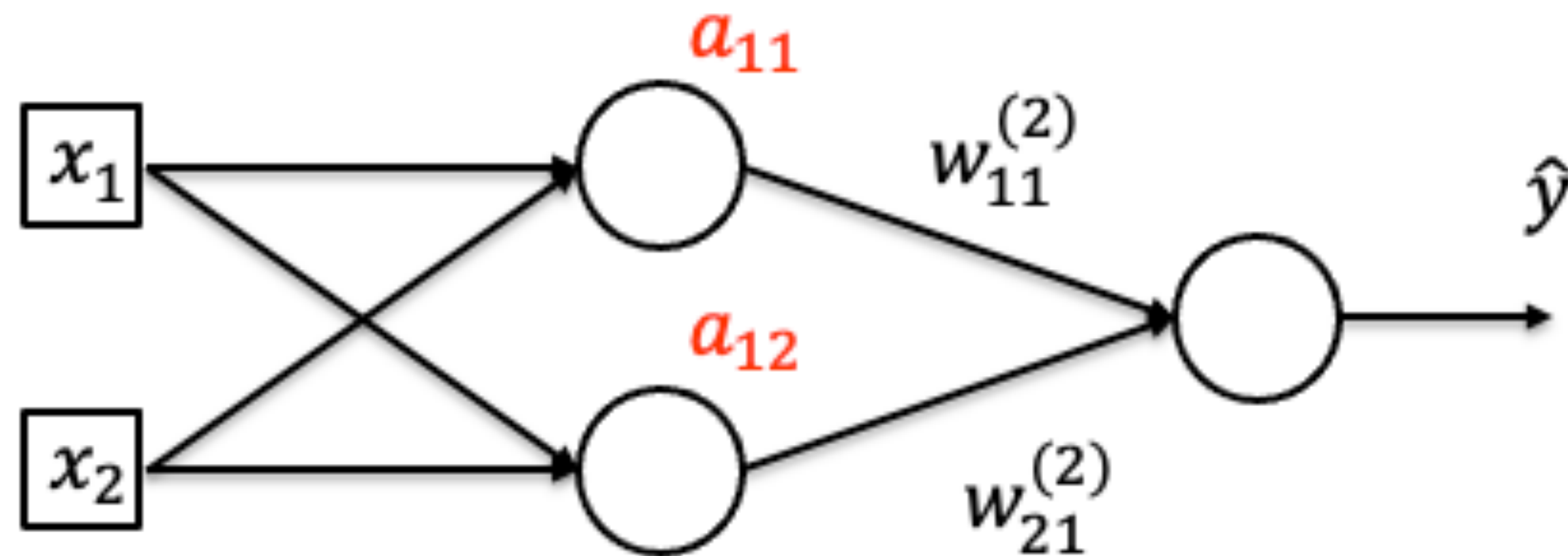
- By chain rule:  $\dfrac{\partial l}{\partial w_{11}} = (\hat{y} - y)x_1$

# Calculate Gradient (on one data point)



$$\frac{\partial \hat{y}}{\partial z} = \sigma'(z) = \sigma(z)(1 - \sigma(z))$$

- By chain rule: $\dfrac{\partial l}{\partial x_1} = \dfrac{\partial l}{\partial \hat{y}} \dfrac{\partial \hat{y}}{\partial z} w_{11} = (\hat{y} - y) w_{11}$

# Calculate Gradient (on one data point)



Make it deeper

$$\frac{\partial \hat{y}}{\partial z} = \sigma'(z) = \sigma(z)(1 - \sigma(z))$$

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

# Calculate Gradient (on one data point)



$$\frac{\sigma(z_{11})}{}$$

$$z_{11} \xrightarrow[\quad\frac{\partial a_{11}}{\partial z_{11}} = \sigma'(z_{11})\quad]{\sigma(z_{11})} a_{11} \xrightarrow[\quad\frac{\partial l}{\partial a_{11}} = (\hat{y} - y)w_{11}^{(2)}\quad]{} l(x, y)$$
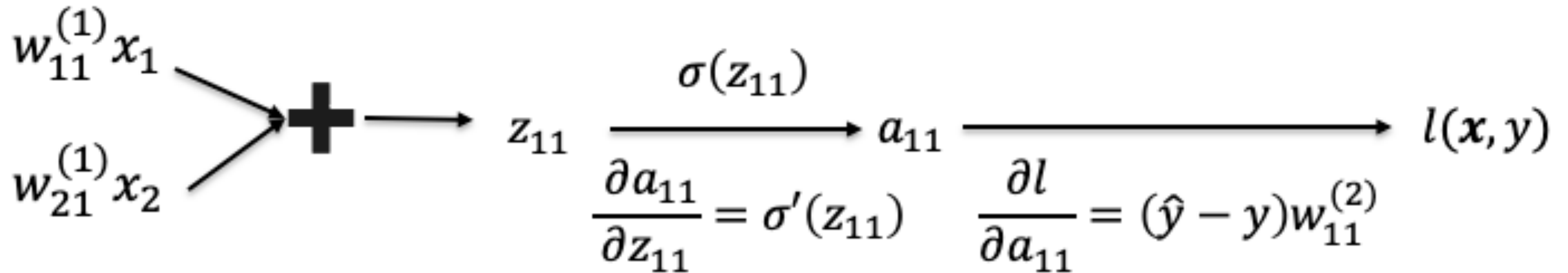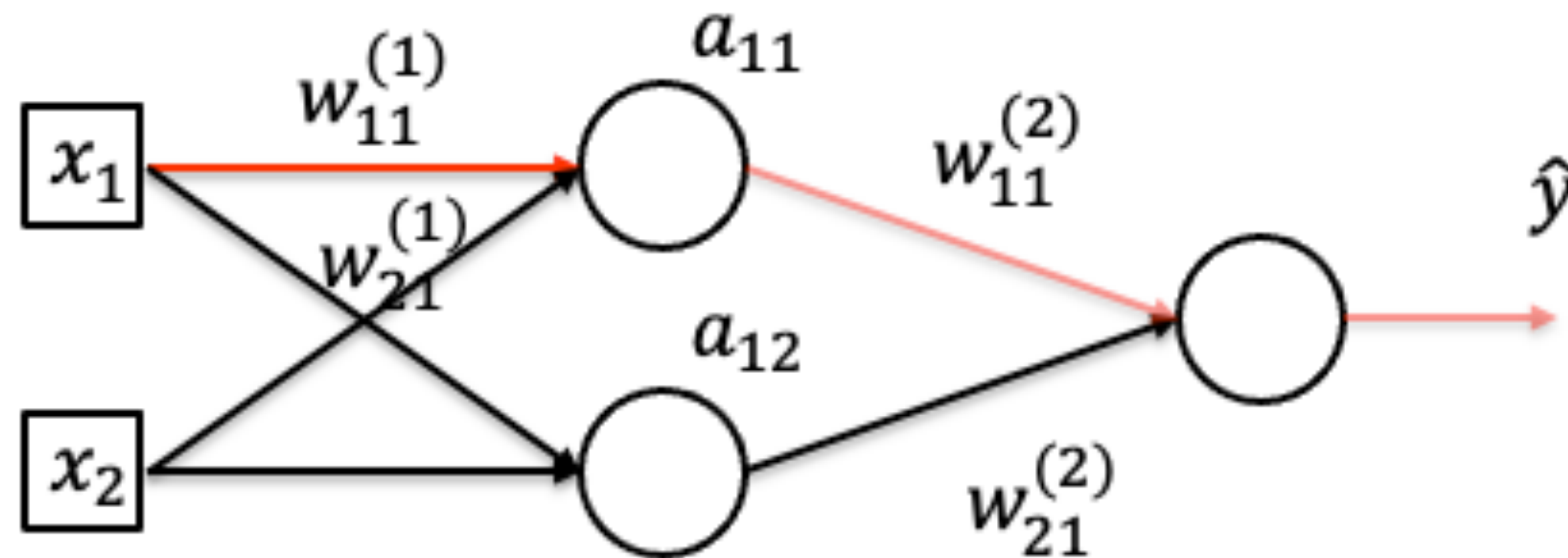
- By chain rule: $\dfrac{\partial l}{\partial w_{11}^{(1)}} = \dfrac{\partial l}{\partial a_{11}} \dfrac{\partial a_{11}}{\partial w_{11}^{(1)}} = (\hat{y} - y)w_{11}^{(2)} \dfrac{\partial a_{11}}{\partial w_{11}^{(1)}}$

# Calculate Gradient (on one data point)



$$w_{11}^{(1)} x_1$$
$$w_{21}^{(1)} x_2$$
$$\boldsymbol{+} \longrightarrow$$

$$z_{11} \xrightarrow{\sigma(z_{11})} a_{11} \xrightarrow{\phantom{xxxxxxxxxxx}} l(x, y)$$

$$\frac{\partial a_{11}}{\partial z_{11}} = \sigma'(z_{11}) \qquad \frac{\partial l}{\partial a_{11}} = (\hat{y} - y) w_{11}^{(2)}$$

- By chain rule: $\dfrac{\partial l}{\partial w_{11}^{(1)}} = \dfrac{\partial l}{\partial a_{11}} \dfrac{\partial a_{11}}{\partial w_{11}^{(1)}} = (\hat{y} - y) w_{11}^{(2)} a_{11}(1 - a_{11}) x_1$
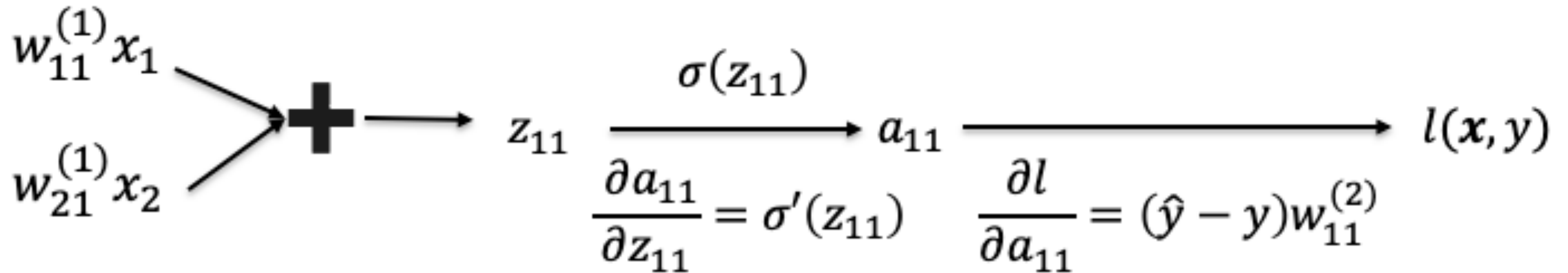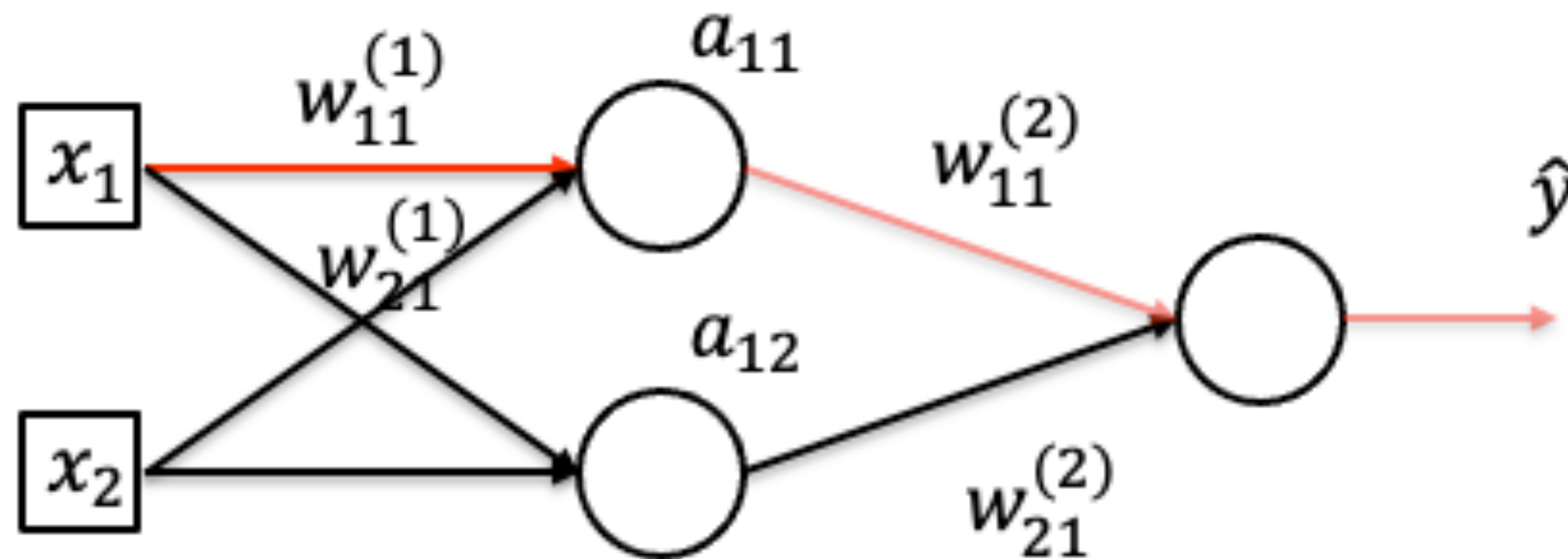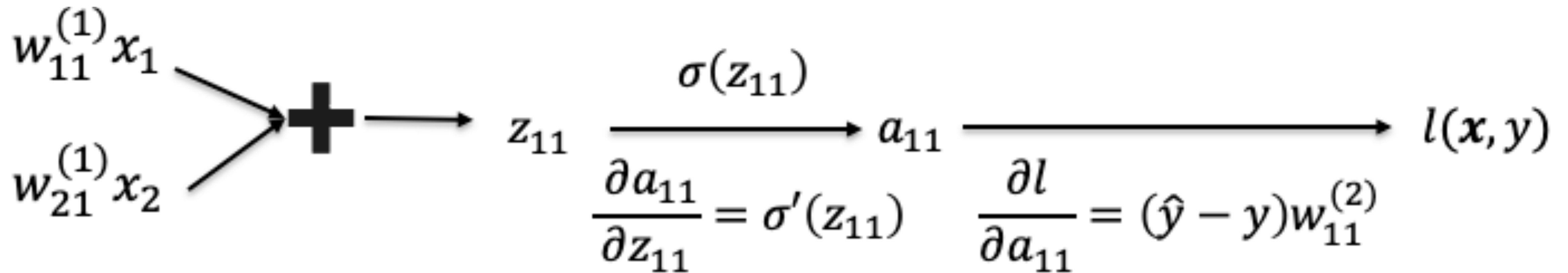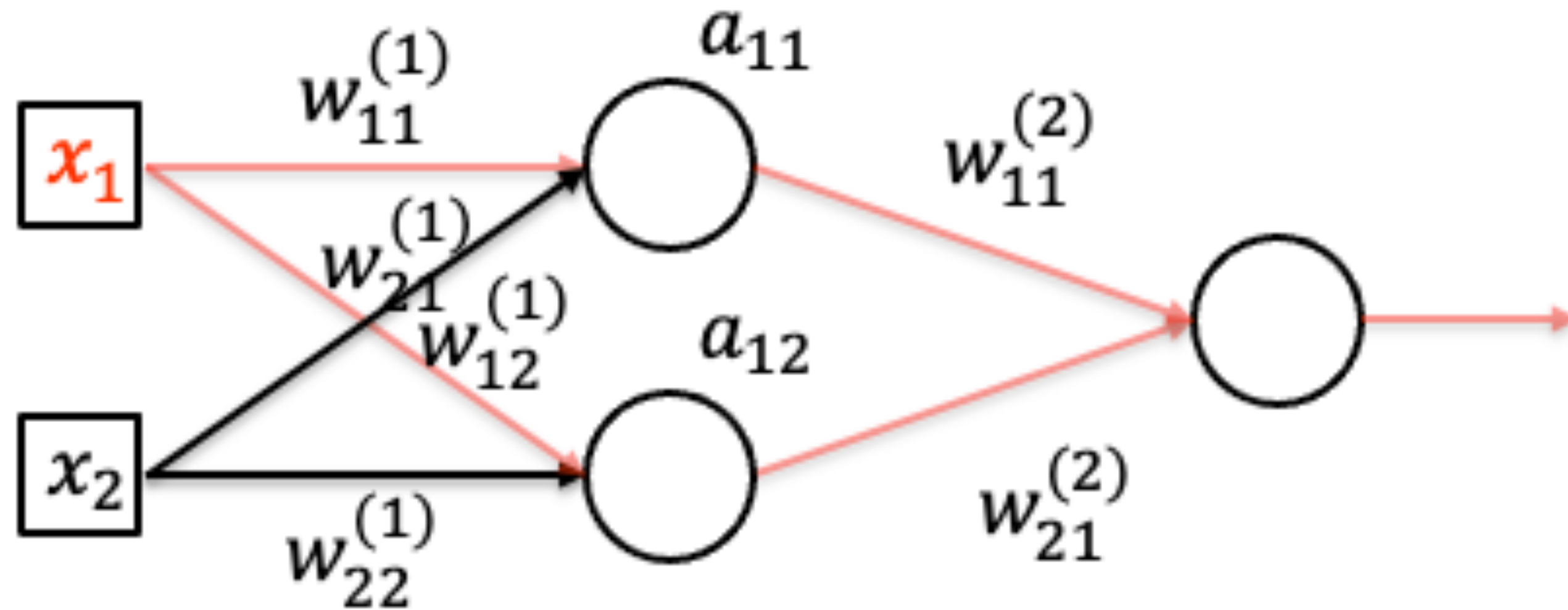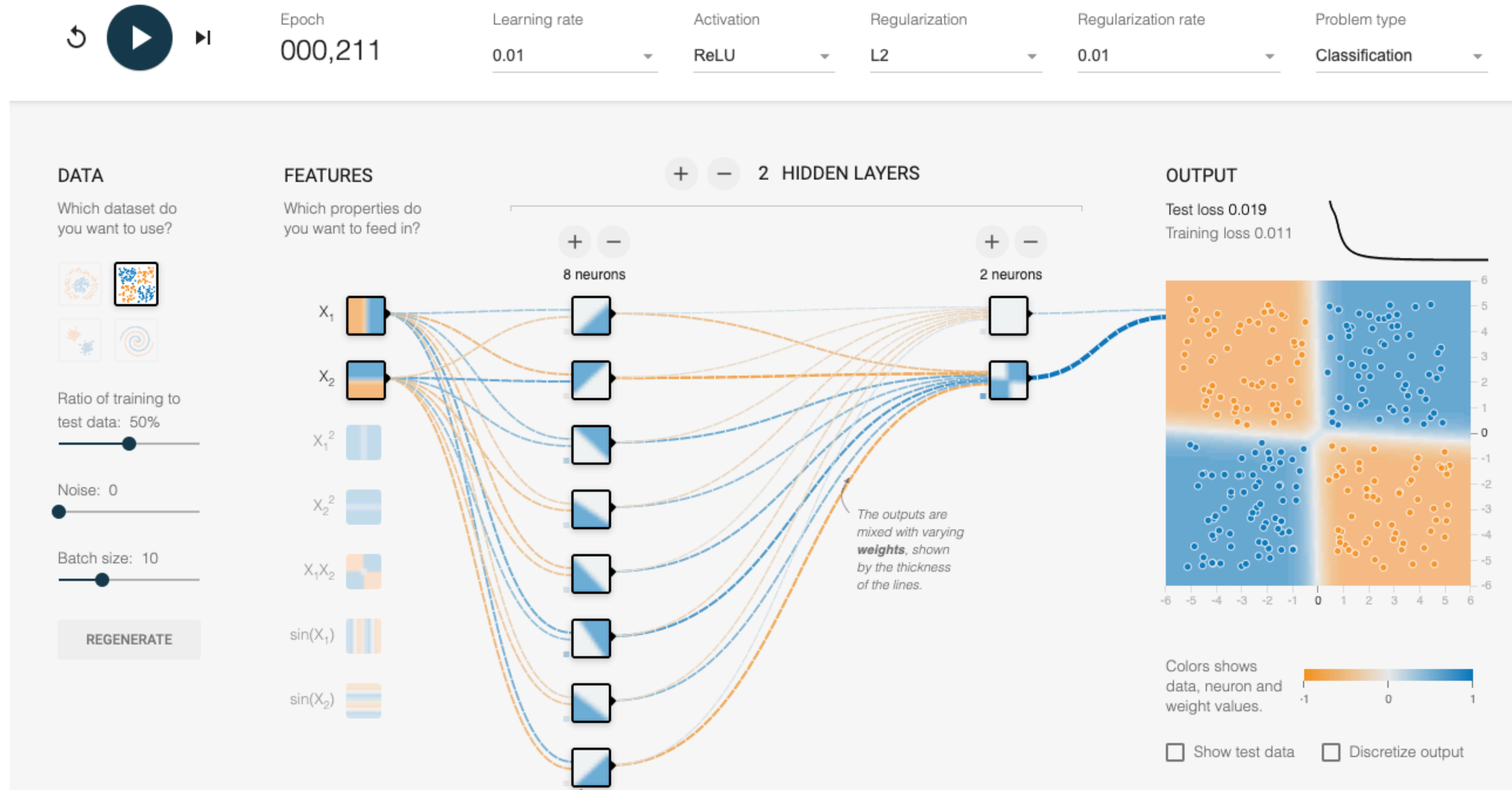
# Calculate Gradient (on one data point)



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

In the diagram:

$x_1$ — $w_{11}^{(1)}$ — $a_{11}$ — $w_{11}^{(2)}$

$w_{21}^{(1)}$

$w_{12}^{(1)}$ — $a_{12}$

$x_2$ — $w_{22}^{(1)}$ — $w_{21}^{(2)}$

$w_{11}^{(1)}x_1$

$w_{21}^{(1)}x_2$

$z_{11}$  $\xrightarrow{\sigma(z_{11})}$  $a_{11}$  $\longrightarrow$  $l(x, y)$

$$\frac{\partial a_{11}}{\partial z_{11}} = \sigma'(z_{11})$$

$$\frac{\partial l}{\partial a_{11}} = (\hat{y} - y)w_{11}^{(2)}$$

- By chain rule:

# 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