



CS 540 Introduction to Artificial Intelligence

Neural Networks (II)

University of Wisconsin-Madison
Spring 2026 Sections 1 & 2

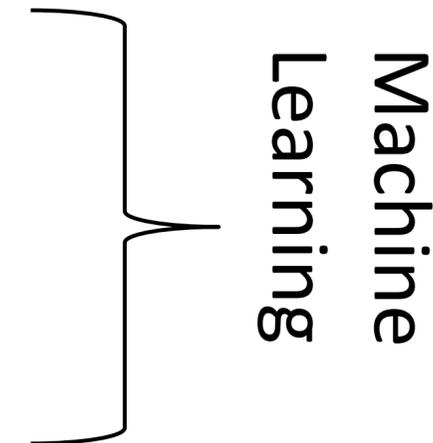
Announcements

- **Homework:**

- HW5 online, due **Wednesday March 4th at 11:59 PM**
- HW6 will also be released on Wednesday

- **Class roadmap:**

Neural Networks II
Neural Networks III
Deep Learning I



Midterm Information

- **Time: March 24th 5:45-7:15 PM**
- **Location (by section **):**
 - Section 001 (Tuesday/Thursday 11-12:15PM): 6210 Social Sciences Bldg
 - Section002 (Tuesday/Thursday 2:30-3:45PM): B10 Ingraham Hall
 - Students with McBurney accommodations will receive an email with additional information.
 - Students who cannot take the exam on the specified time should contact their instructor if they have not done it yet.
- **Format:** MCQ
- **Cheat sheet:** a handwritten single piece of paper, front and back
- **Calculator:** optional, if it doesn't have an Internet connection
- **Bring:** your WISC ID, pencil (No 2 or softer), your 1-sheet notes.
- **Past exam questions:** on Canvas → Files → Past Exams

Today's outline

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

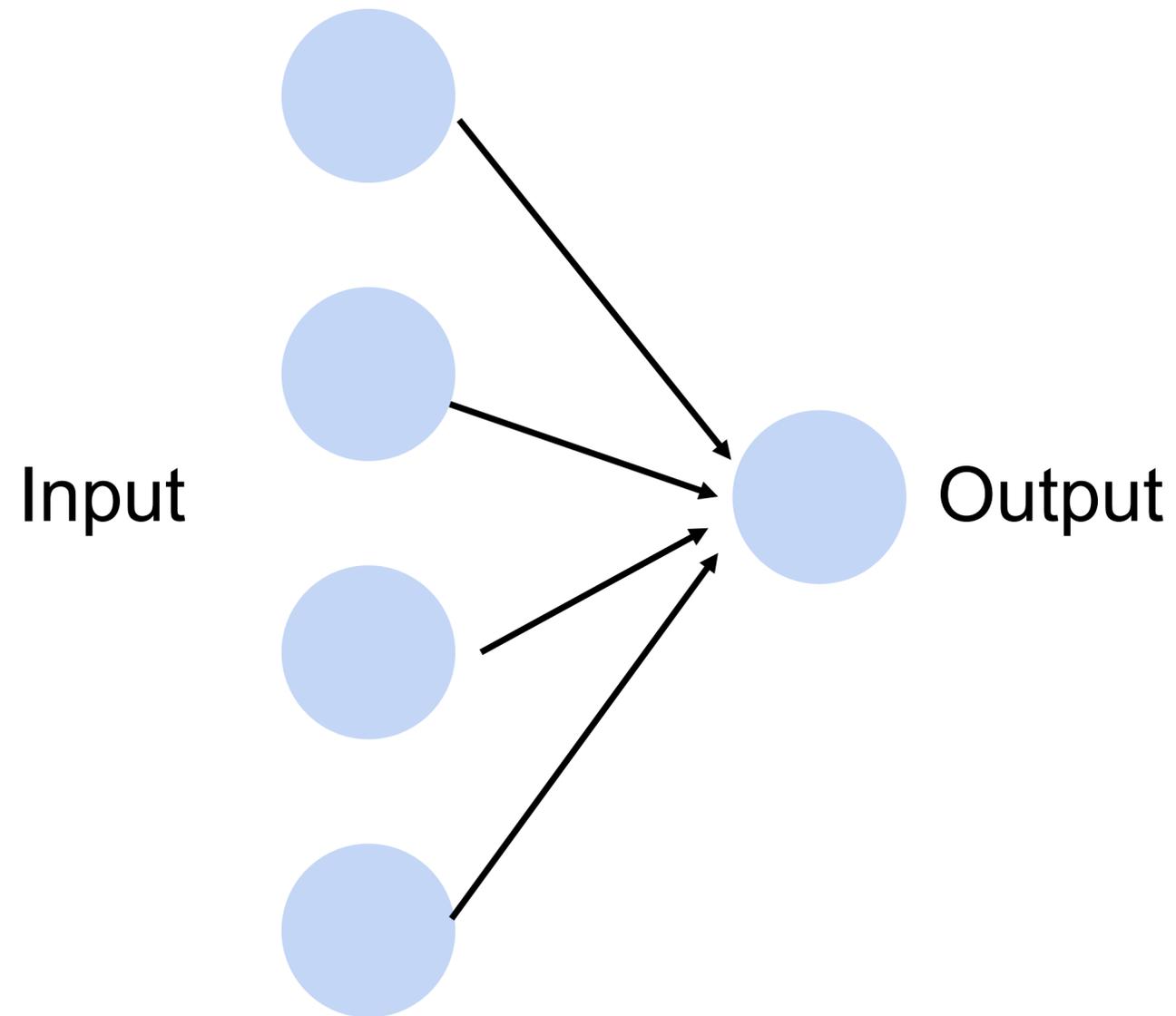
Review: Perceptron

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

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

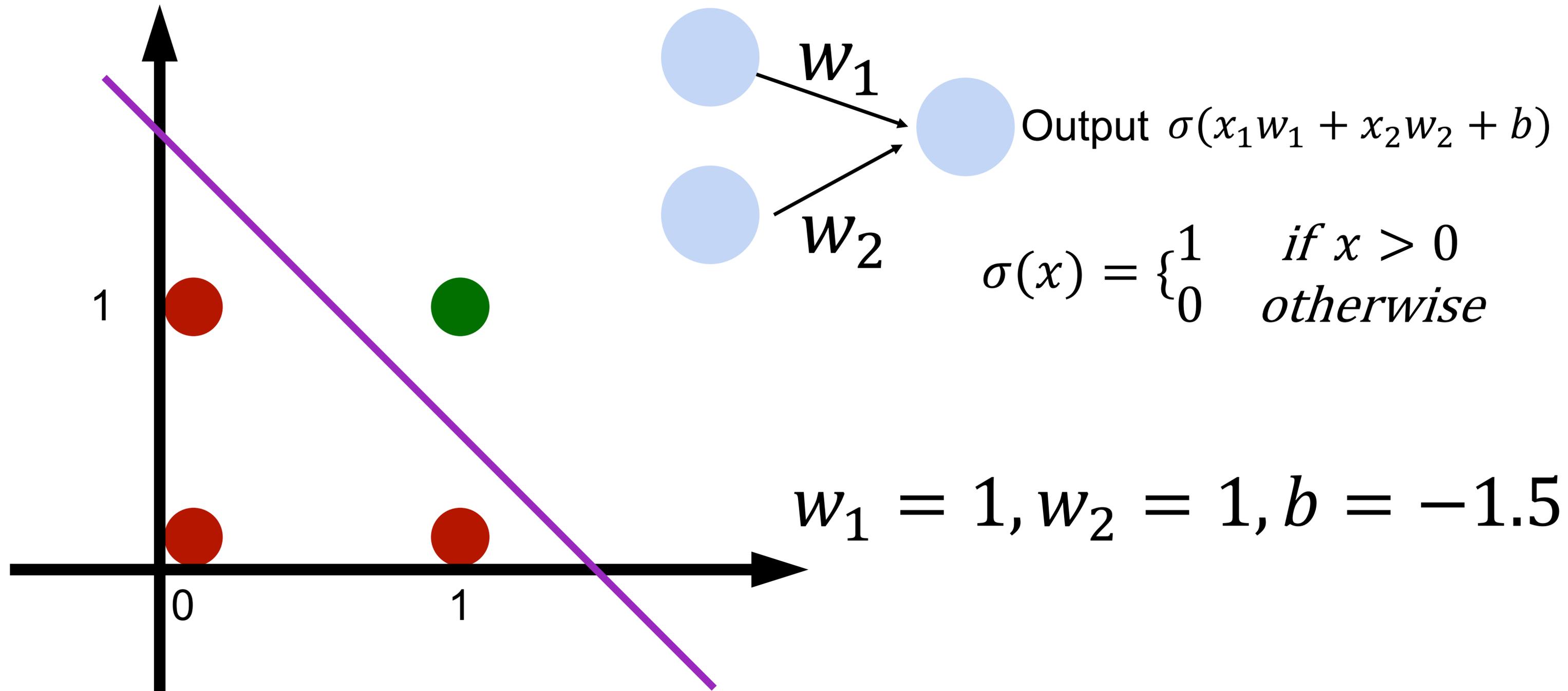
$$\sigma(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases} \text{Activation function}$$

Cats vs. dogs?



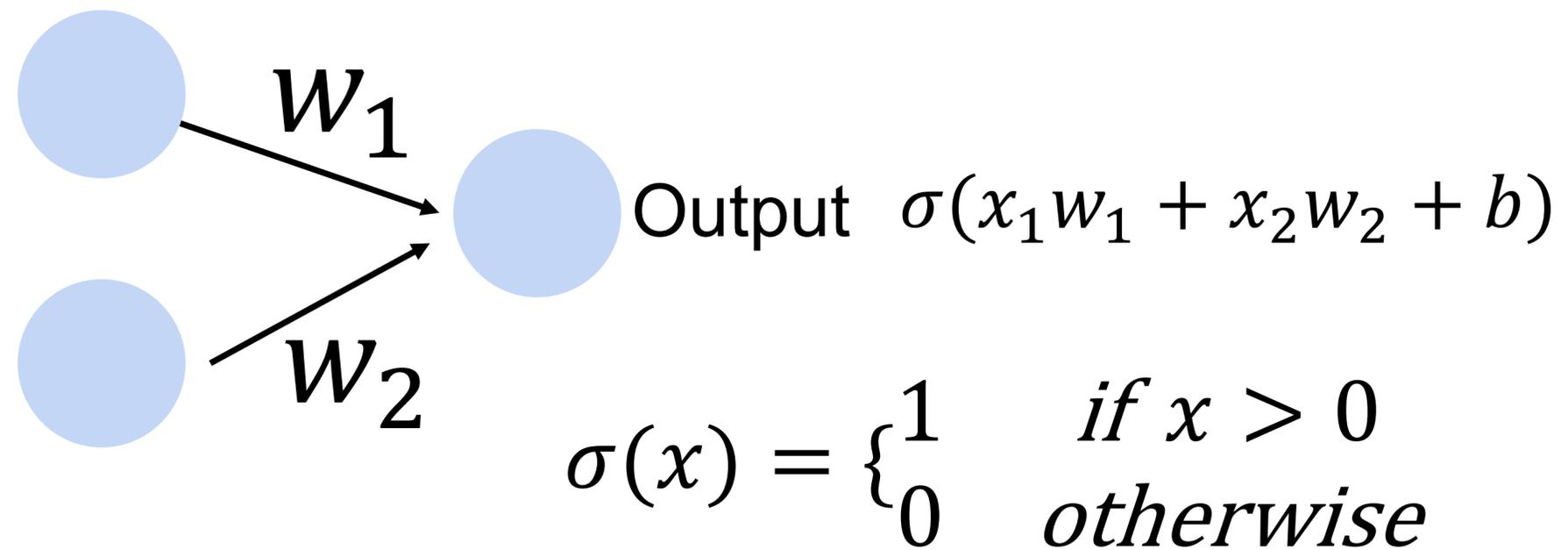
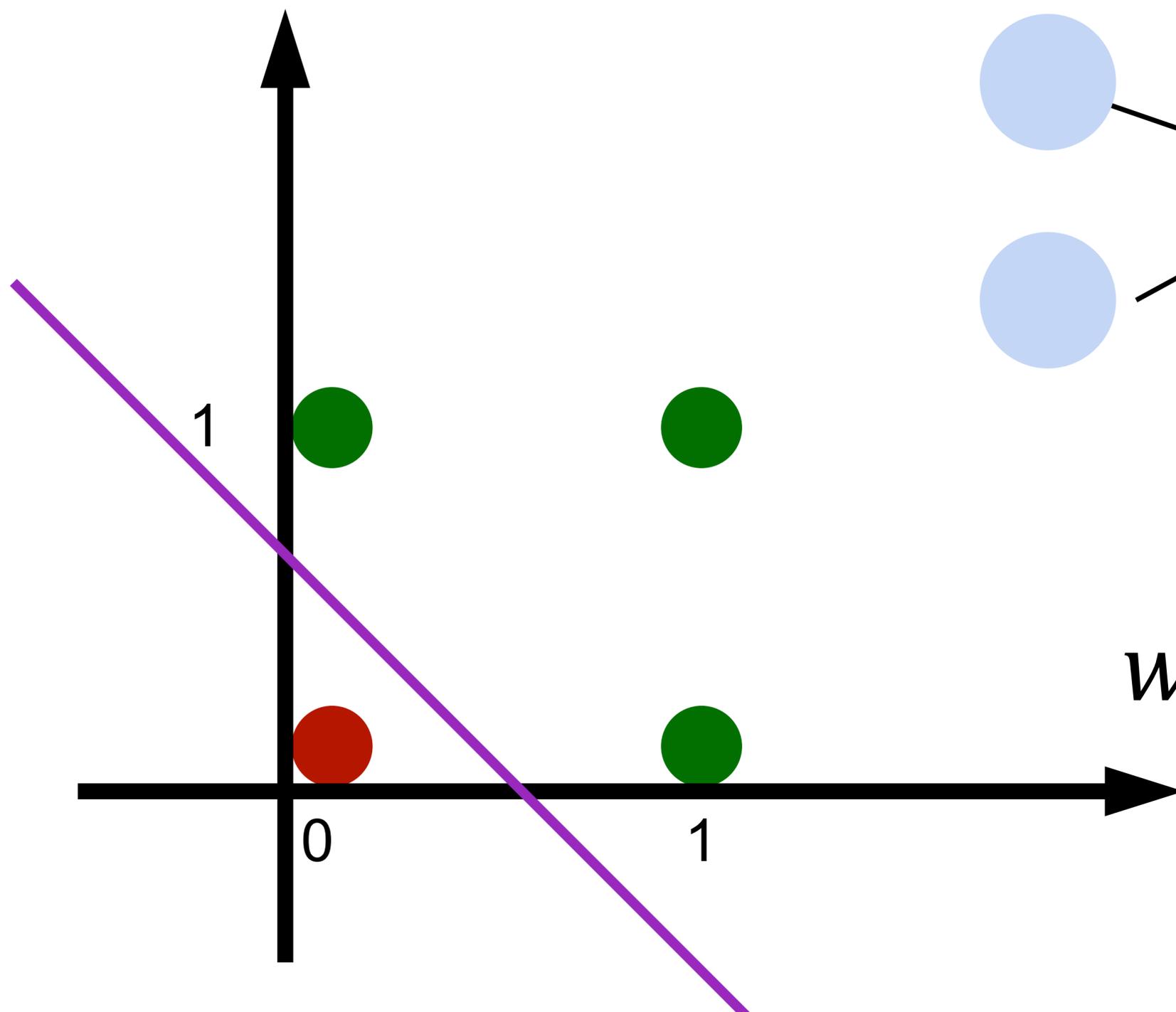
Learning AND function using perceptron

The perceptron can learn an AND function



Learning OR function using perceptron

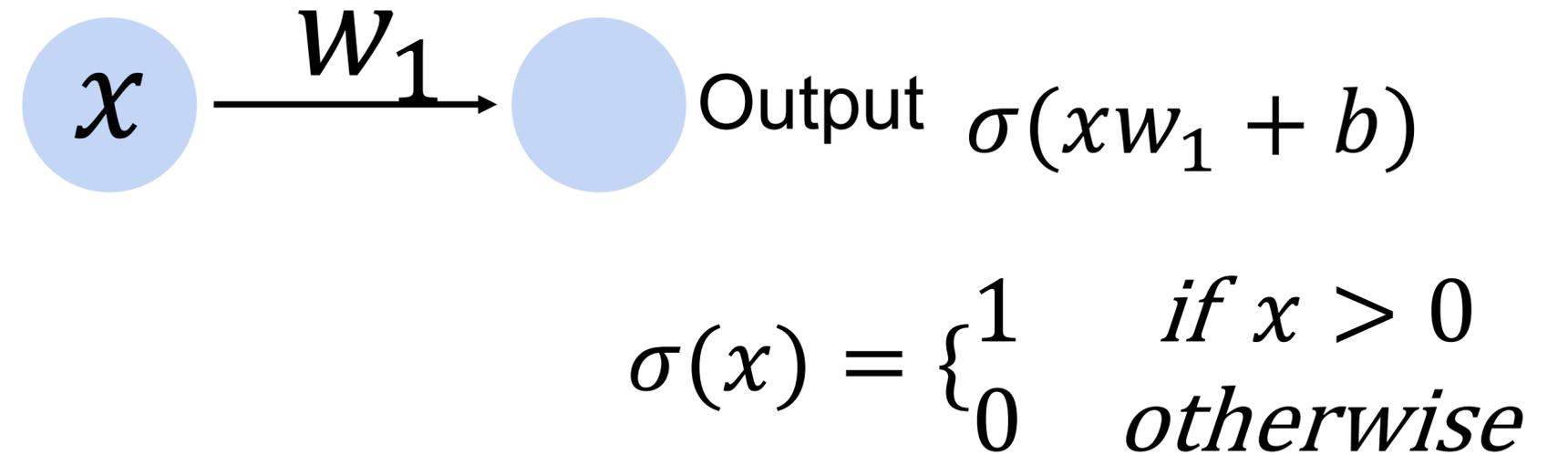
The perceptron can learn an OR function



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

Learning NOT function using perceptron

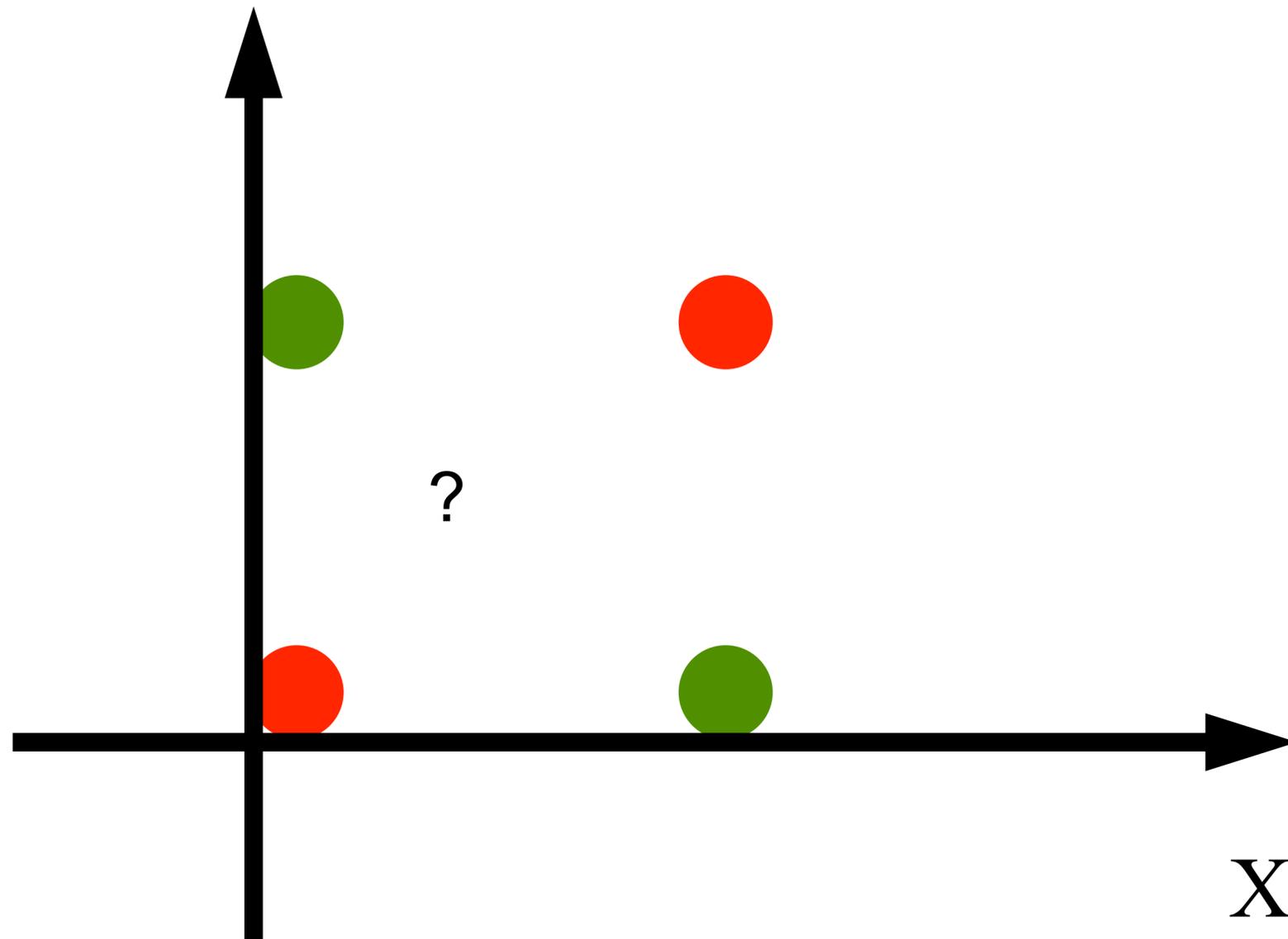
The perceptron can learn NOT function (single input)



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

The limited power of a single neuron

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



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

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

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

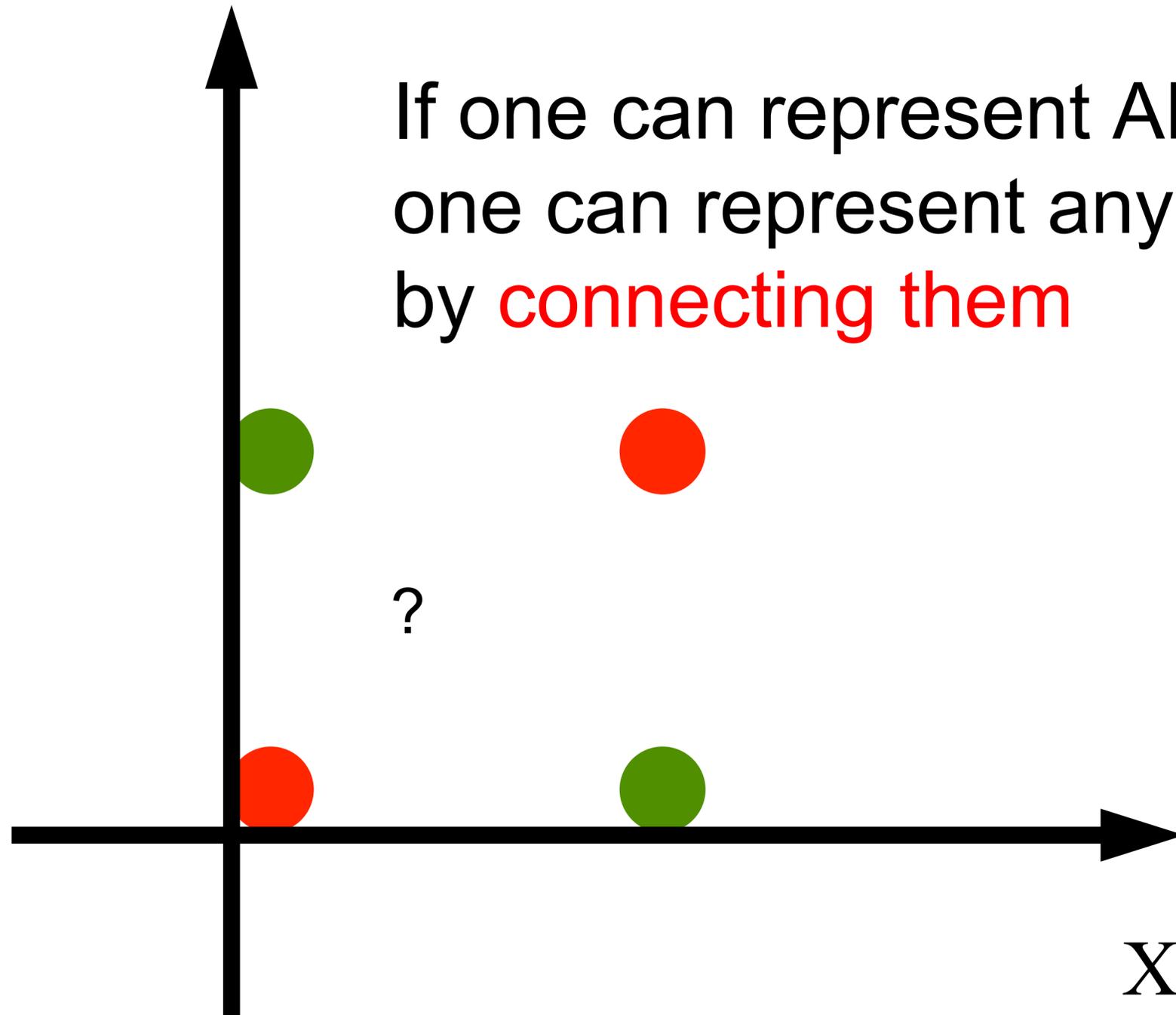
$$x_1 = 0, x_2 = 0, y = 0$$

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

The limited power of a single neuron

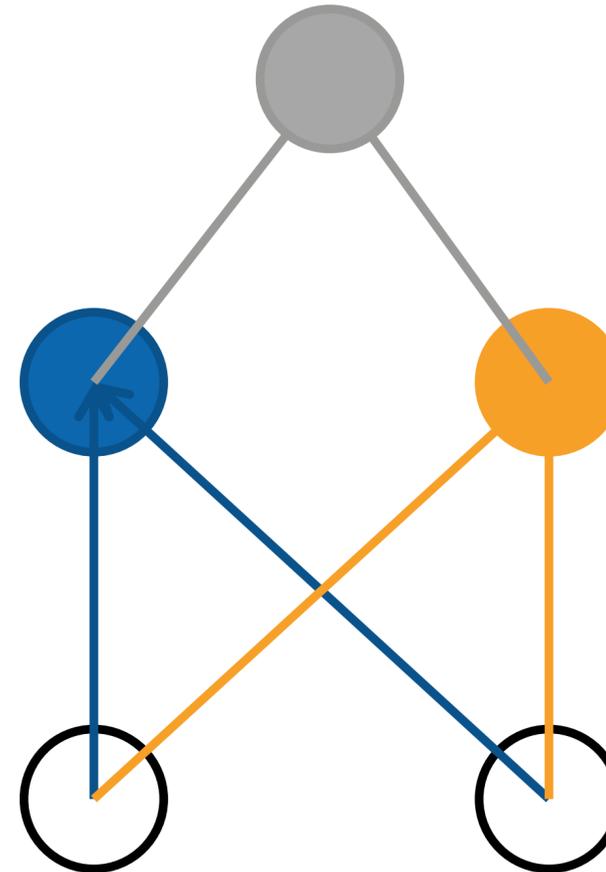
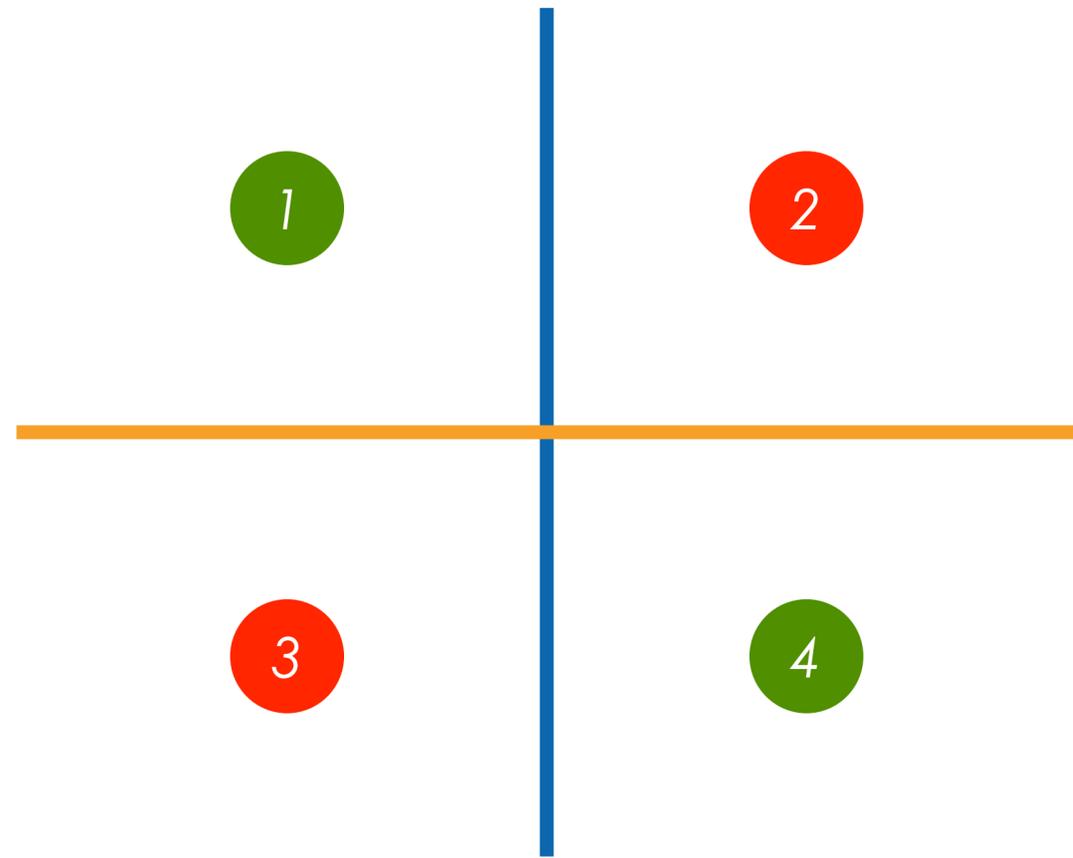
XOR problem

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



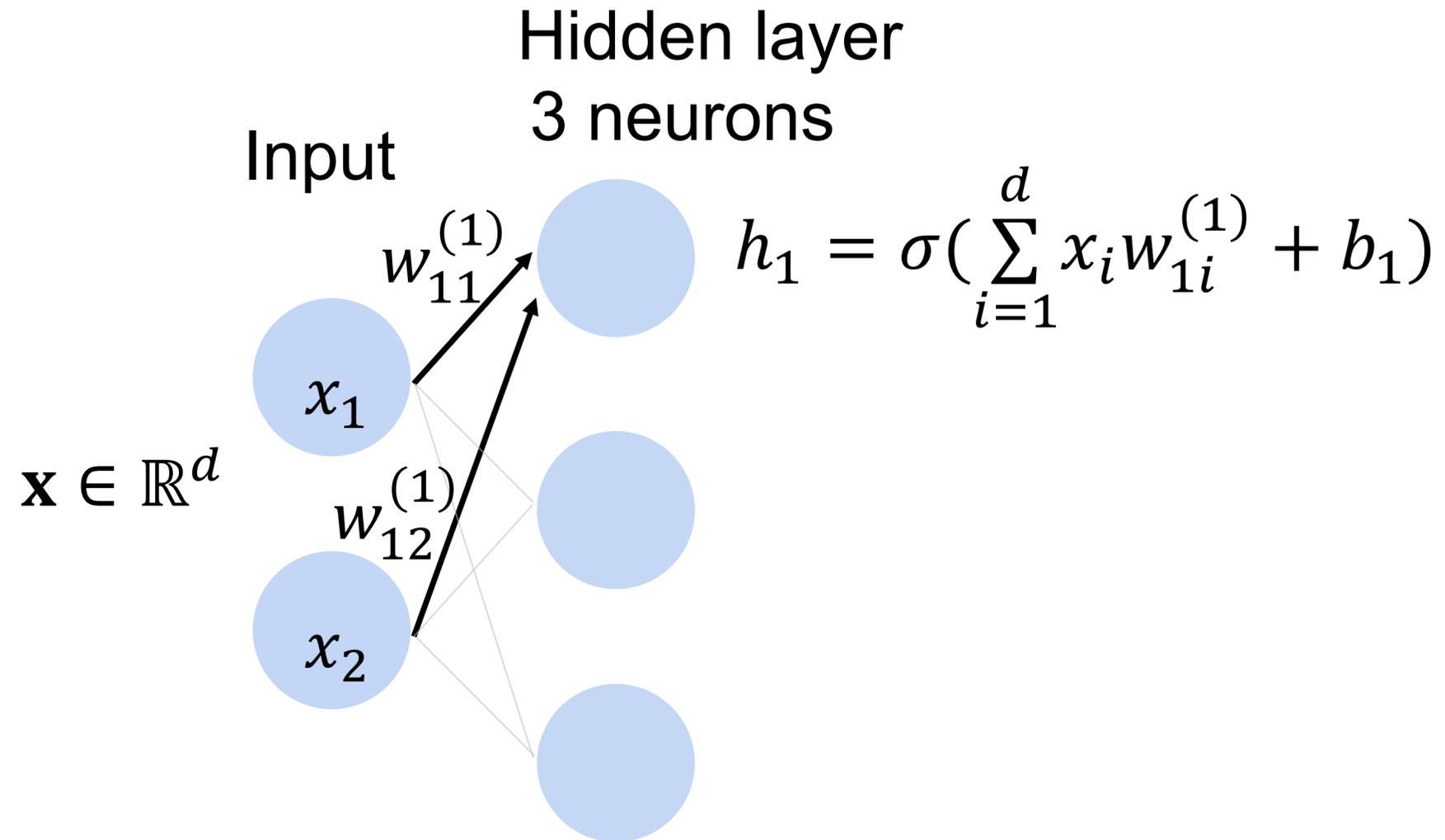
$$\text{XOR}(x_1, x_2) = (x_1 \wedge \neg x_2) \vee (\neg x_1 \wedge x_2)$$

Learning XOR



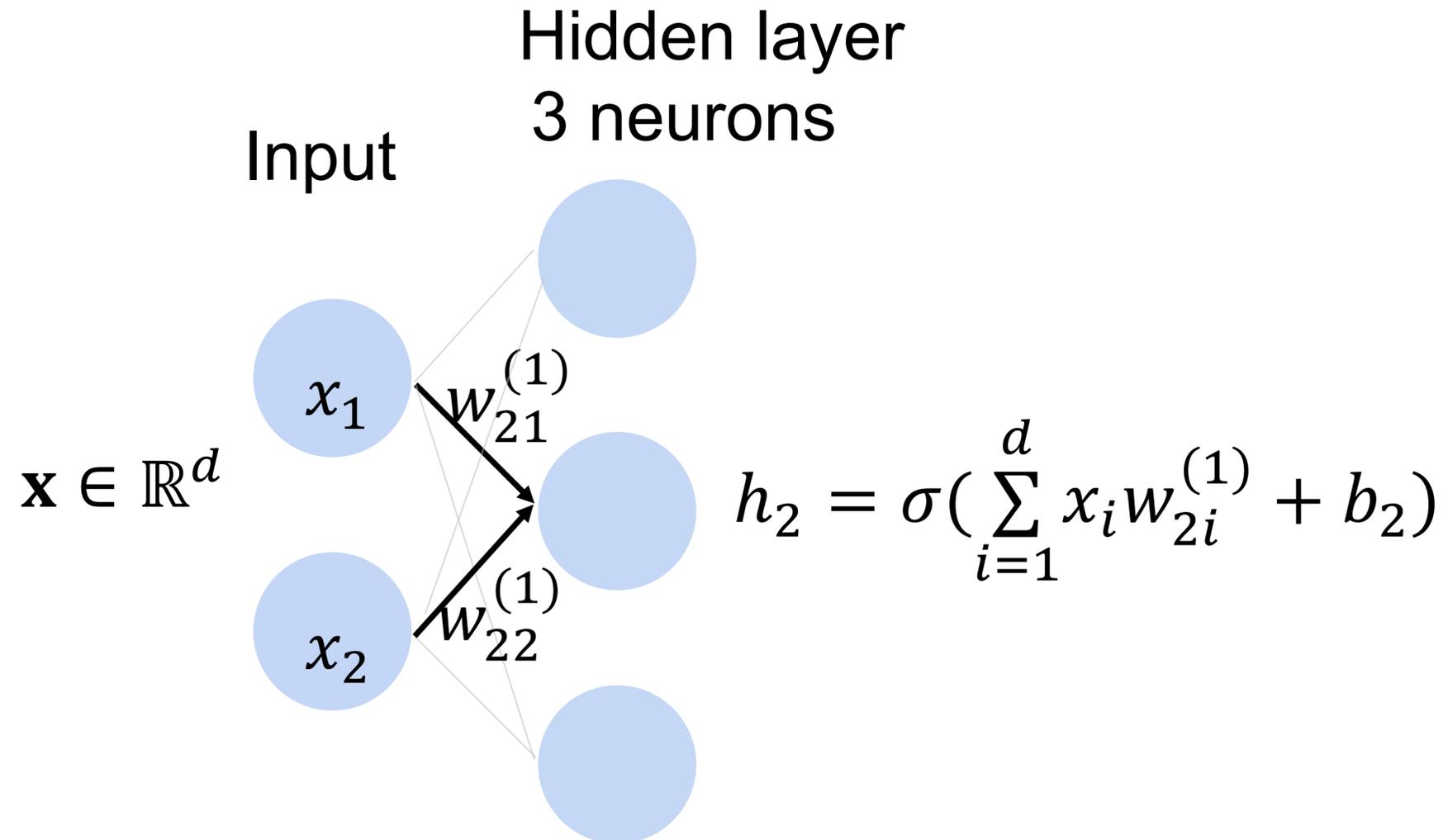
Multi-layer perceptron: Example

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



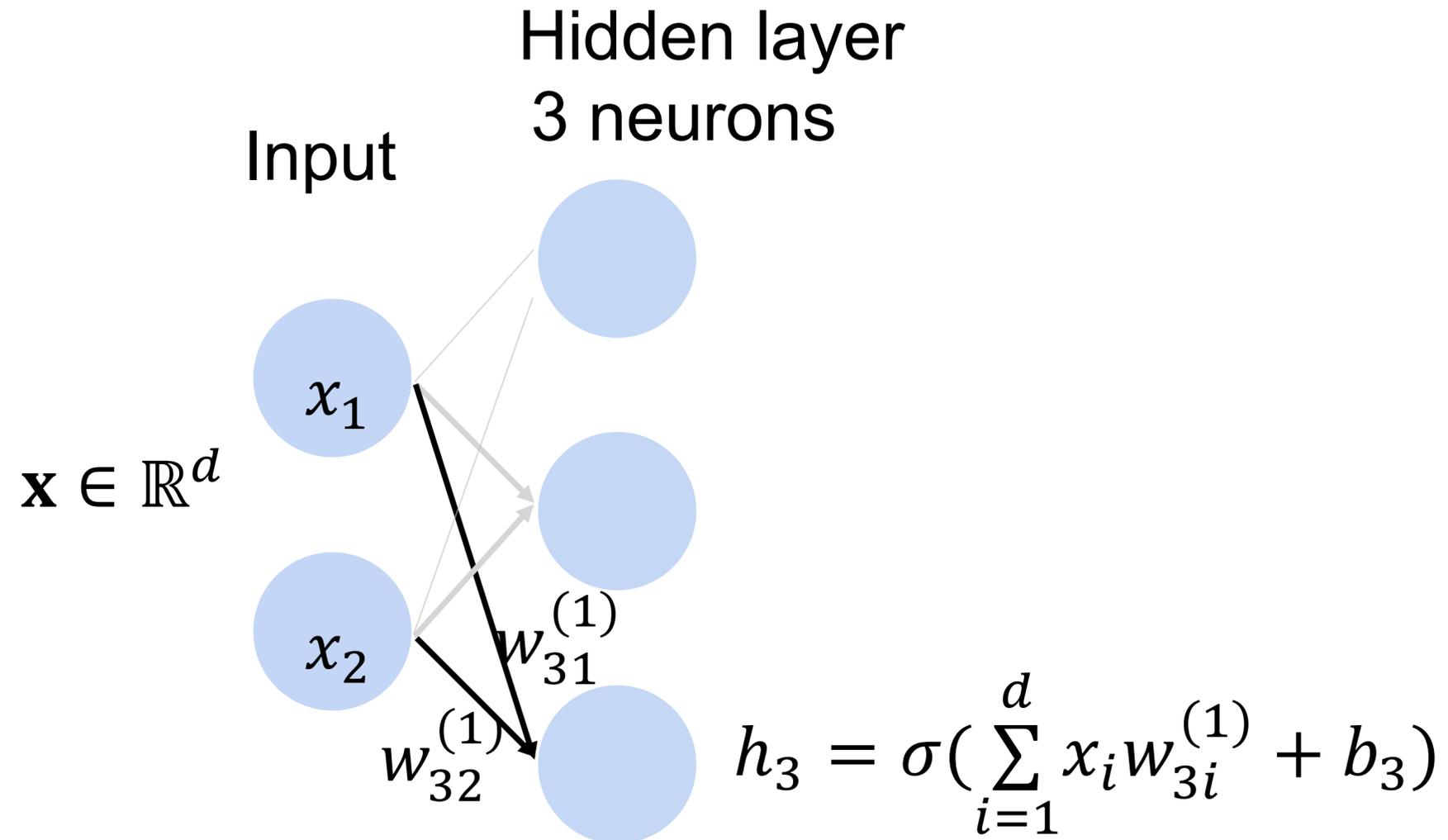
Multi-layer perceptron: Example

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



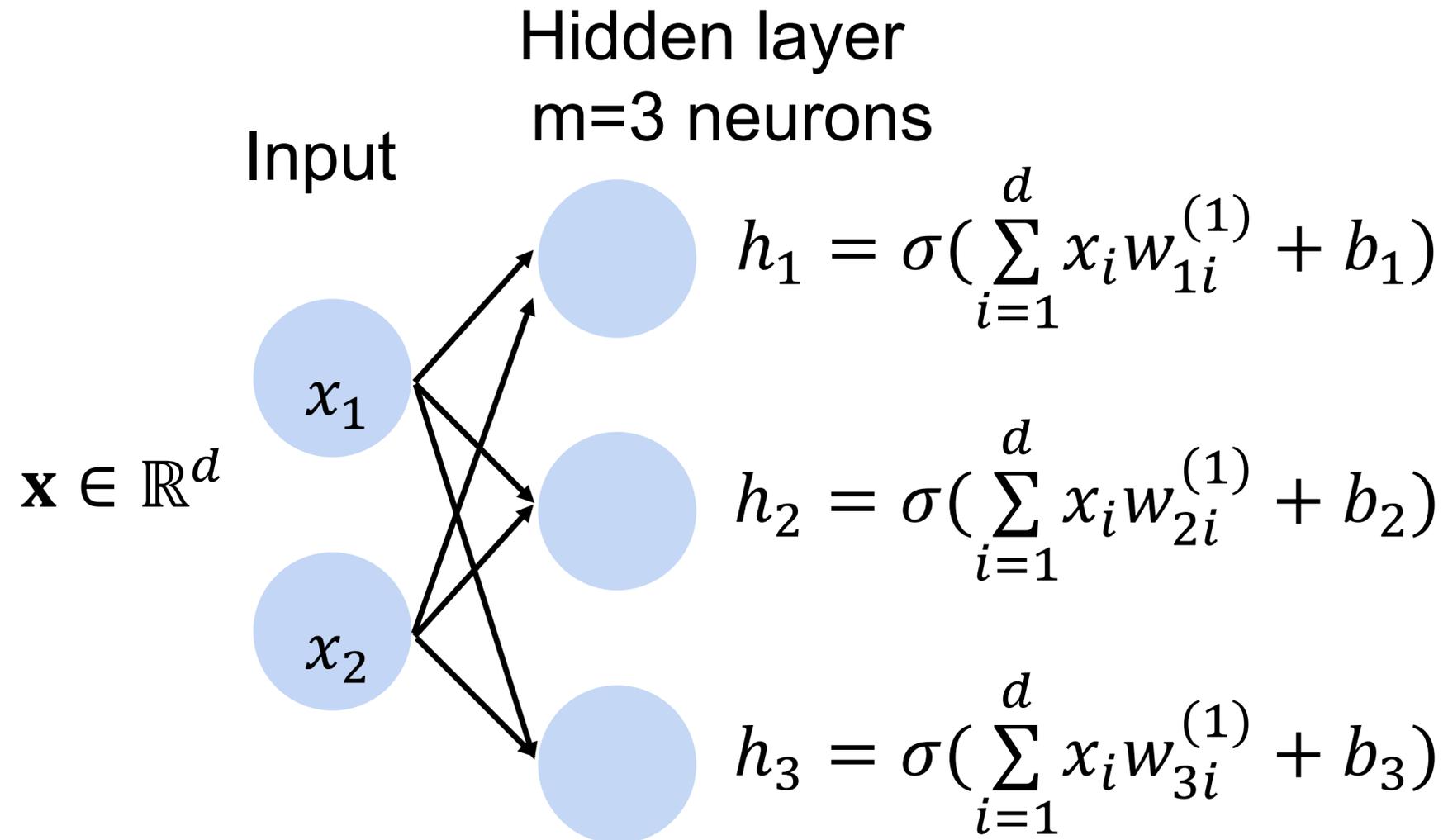
Multi-layer perceptron: Example

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



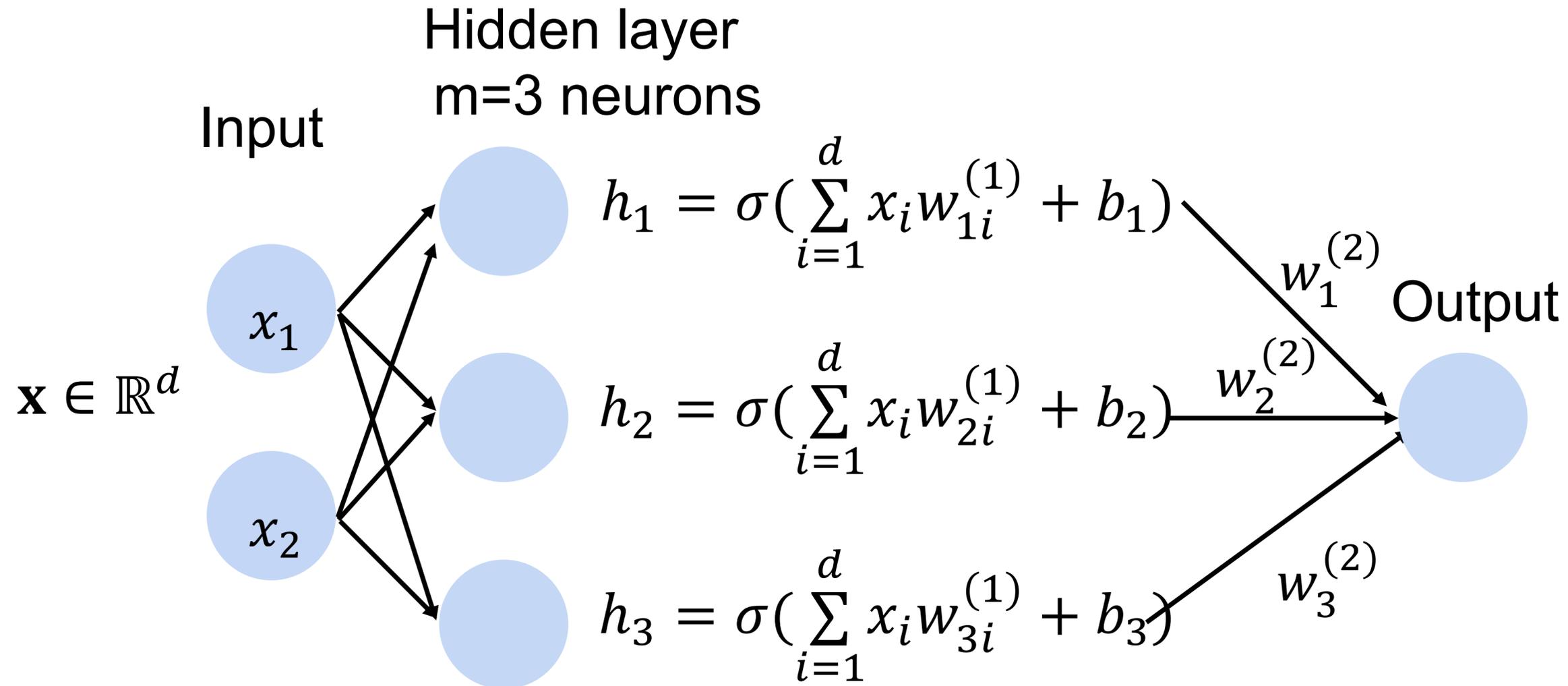
Multi-layer perceptron: Example

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



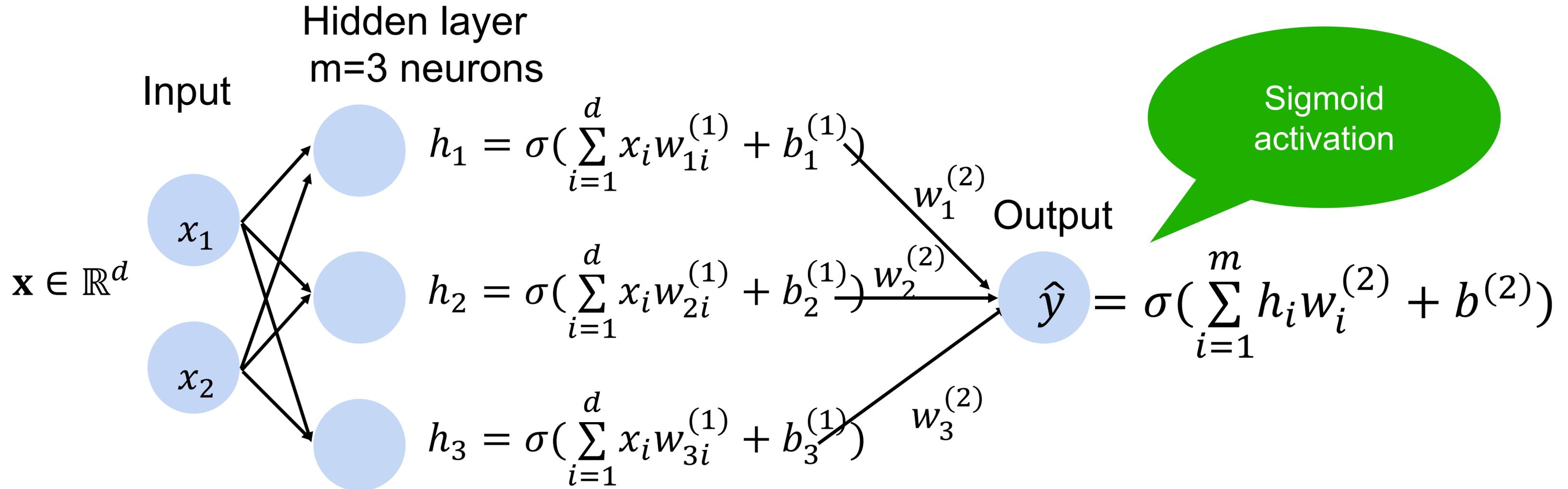
Multi-layer perceptron: Example

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



Multi-layer perceptron: Example

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

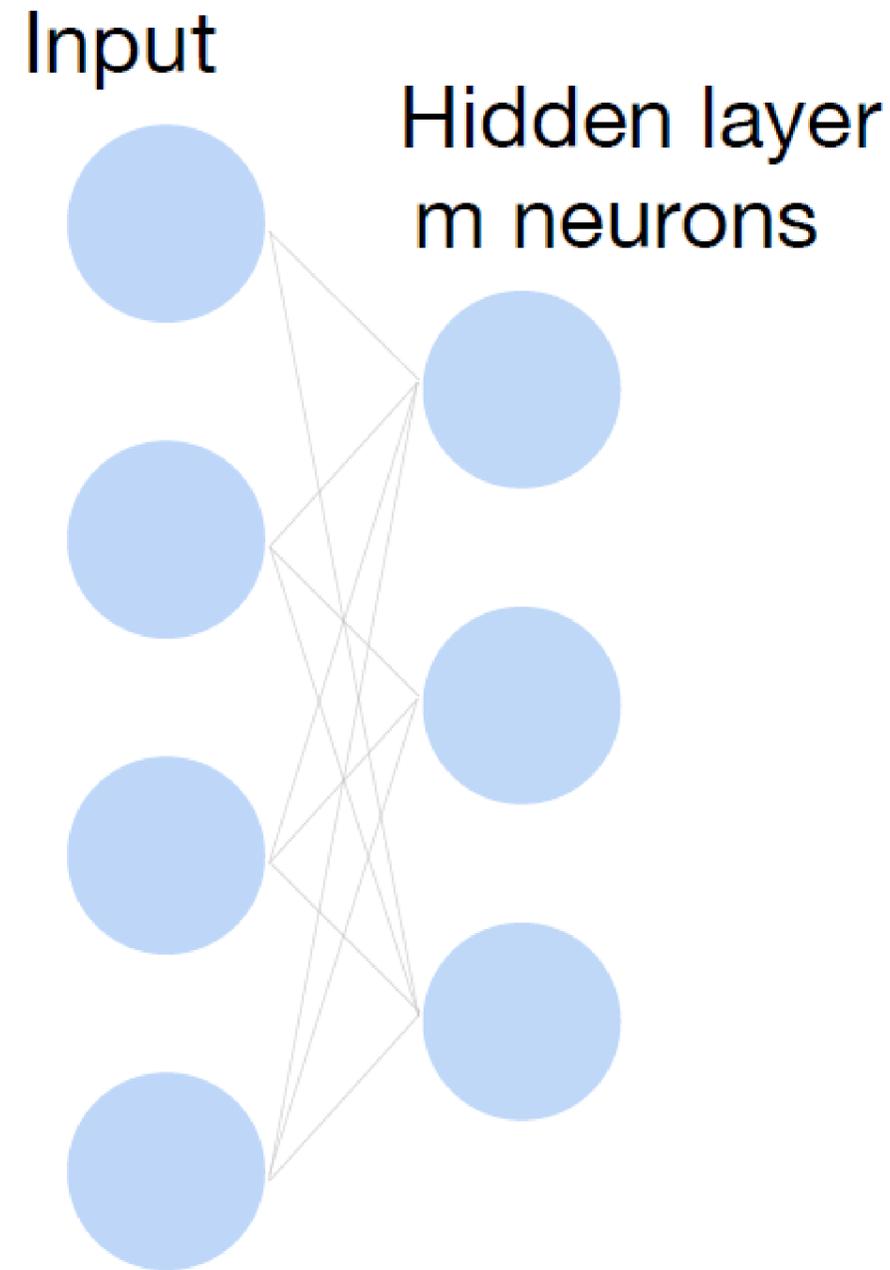


Multi-layer perceptron: Matrix Notation

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

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

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

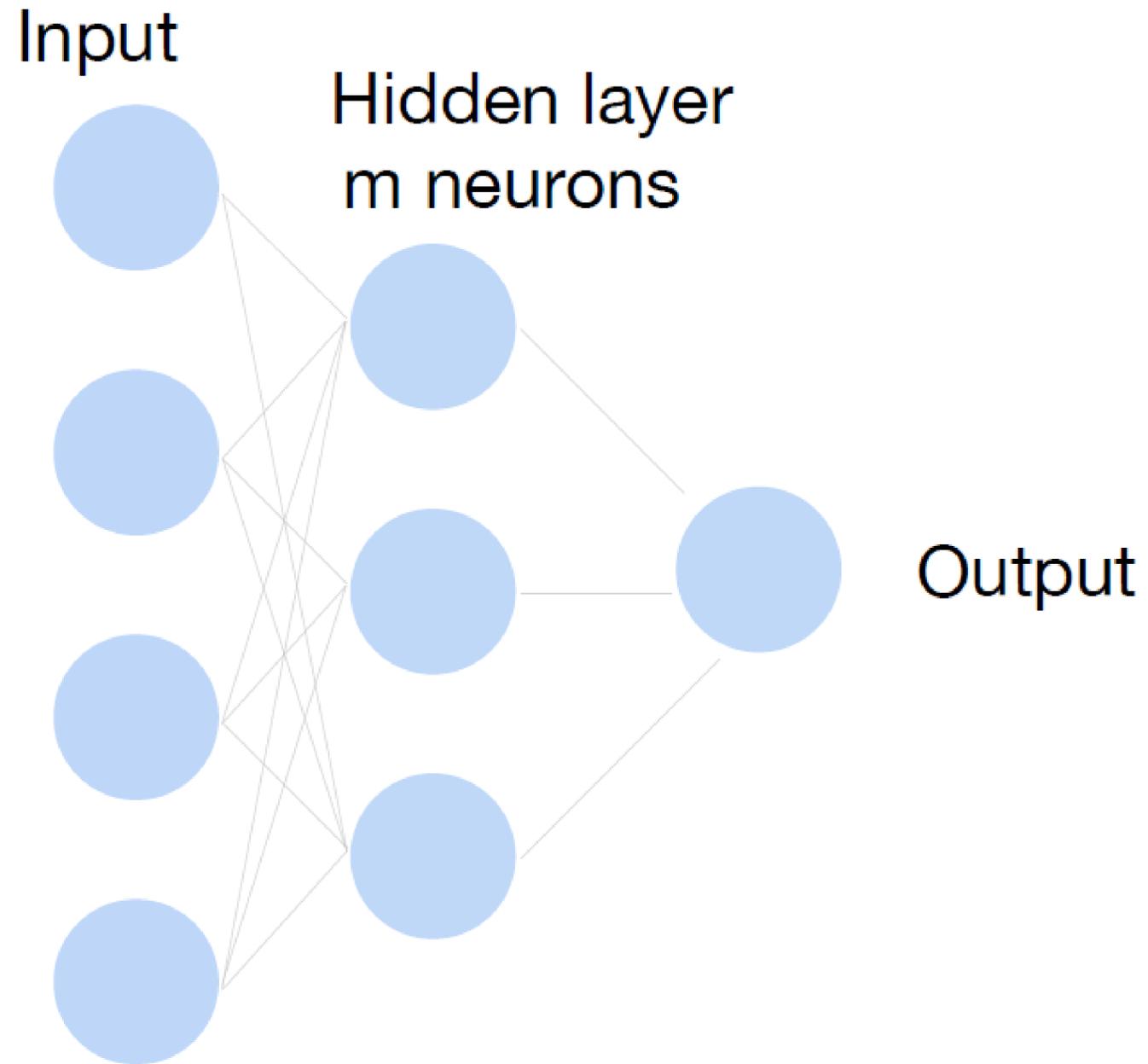


Multi-layer perceptron: Matrix Notation

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

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

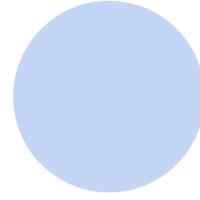
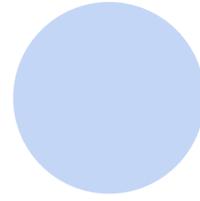
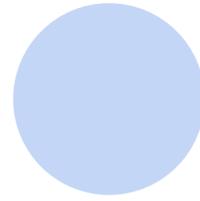
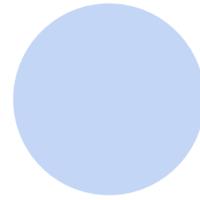
$$\hat{y} = \sigma(\mathbf{w}^{(2)}\mathbf{h} + \mathbf{b}^{(2)})$$



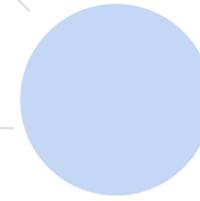
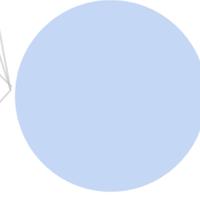
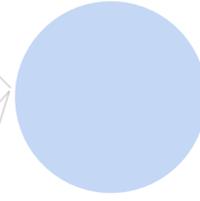
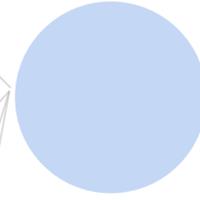
Classify cats vs. dogs



Input

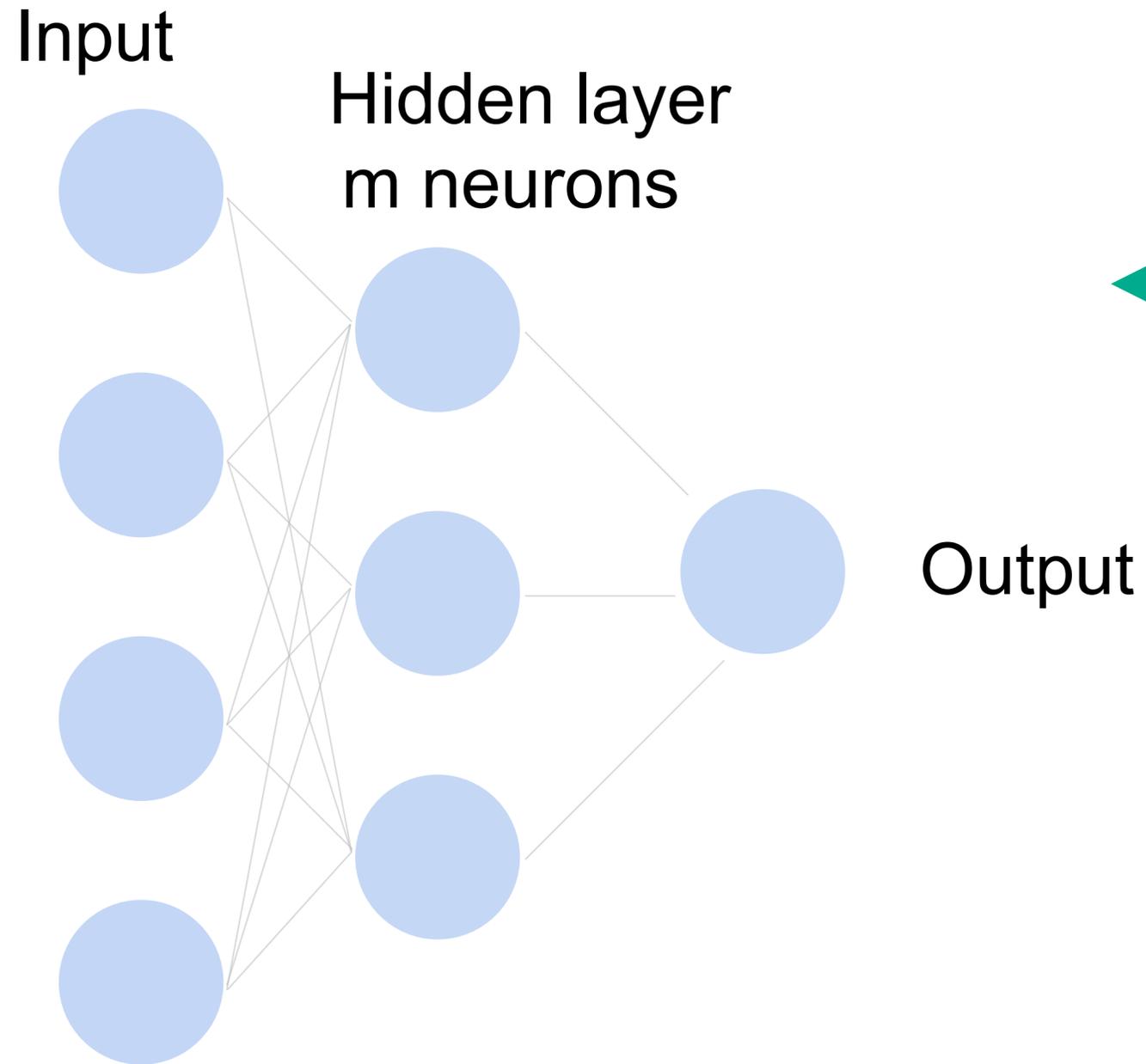


Hidden layer
100 neurons

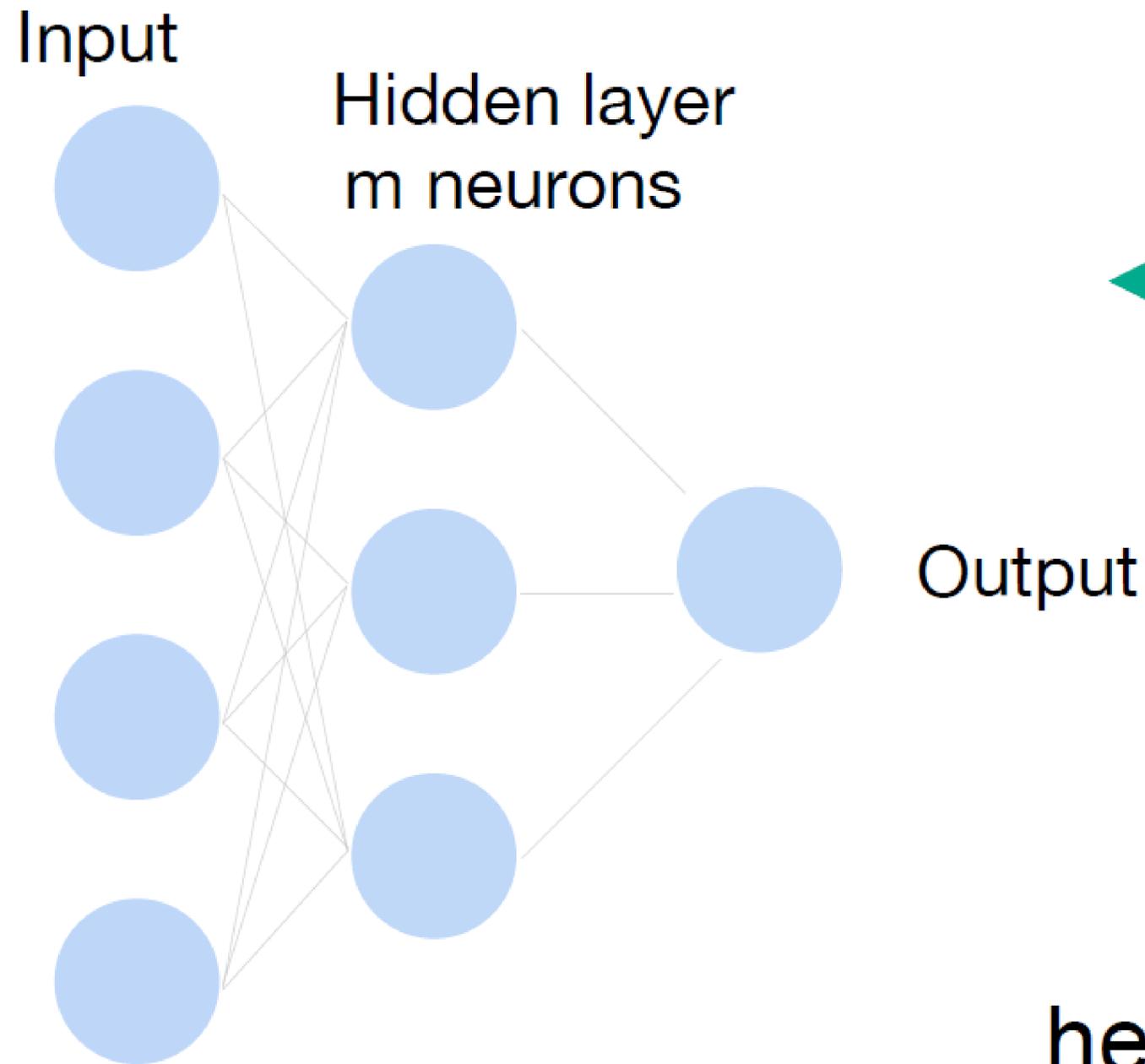


Output

Multi-layer perceptron



Why do we need a nonlinear activation?



Why do we need an a
nonlinear activation?

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

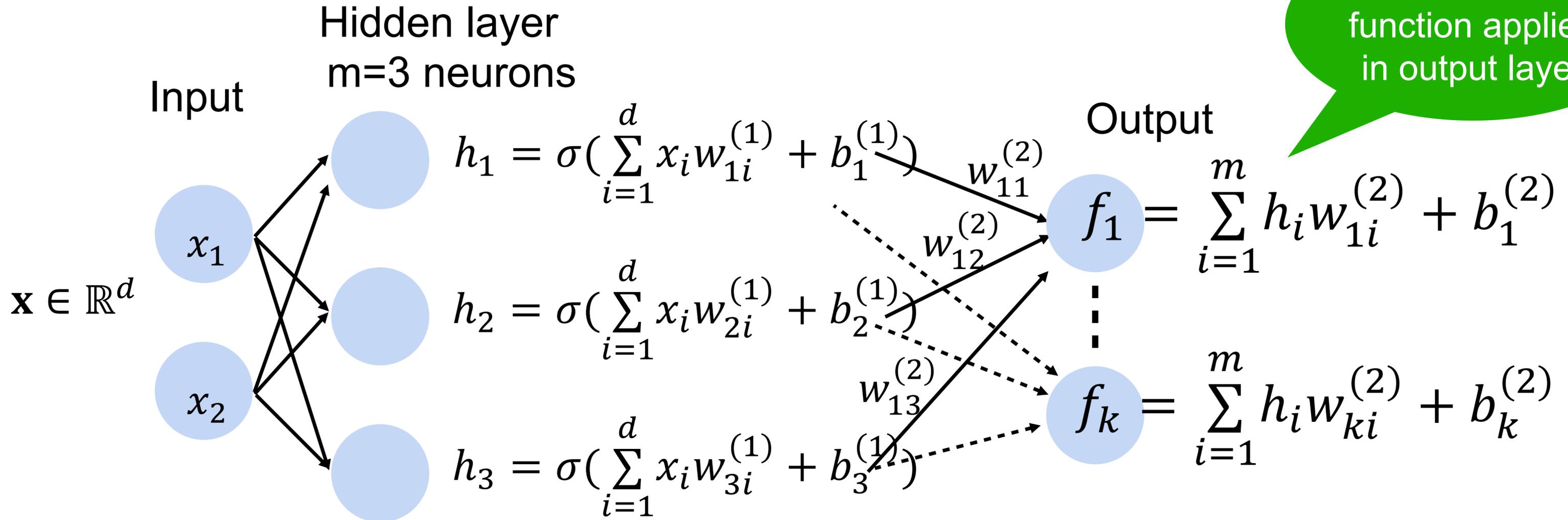
$$f = (\mathbf{w}^{(2)})^T \mathbf{h} + b^{(2)}$$

$$\text{hence } f = (\mathbf{w}^{(2)})^T \mathbf{W}^{(1)}\mathbf{x} + b'$$

Neural network for K-way classification

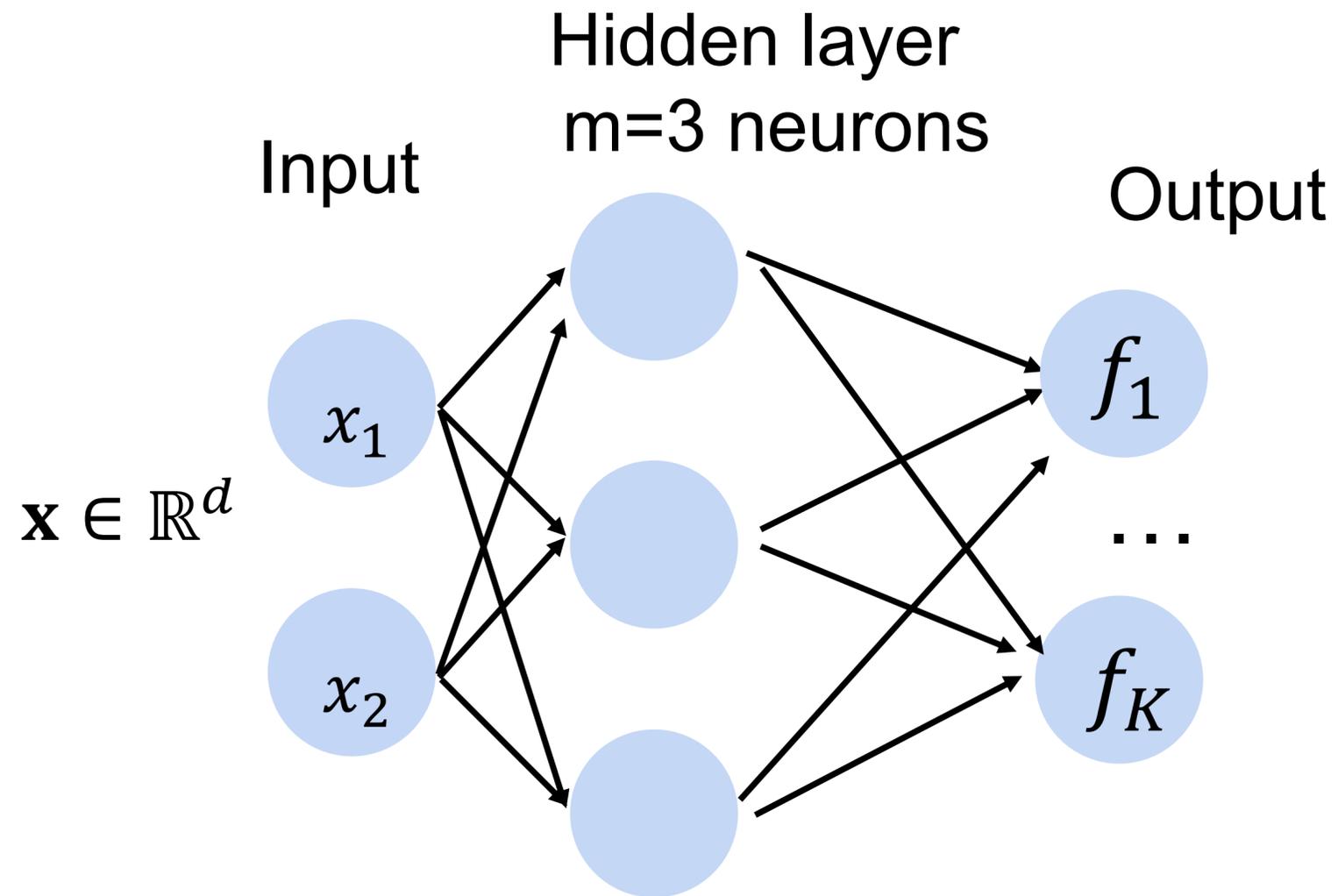
- K outputs in the final layer

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



Softmax

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



$$\begin{aligned} p(y|\mathbf{x}) &= \textit{softmax}(f) \\ &= \frac{\exp(f_y(x))}{\sum_{k=1}^K \exp(f_k(x))} \end{aligned}$$

Softmax

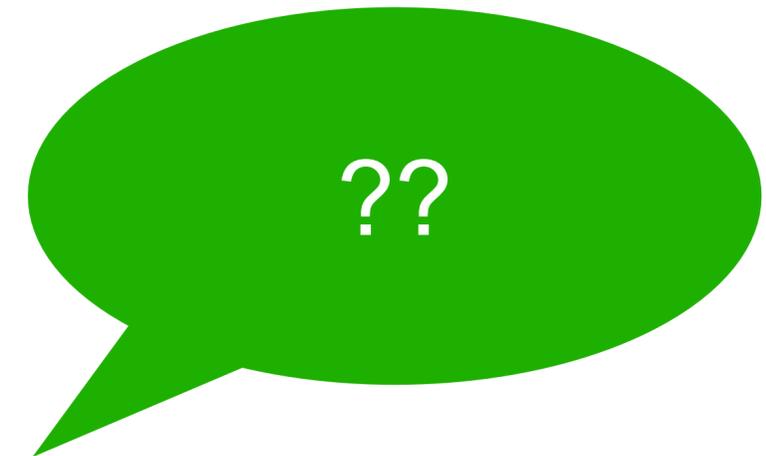
Turns outputs f into probabilities (sum up to 1 across K classes)

Output
layer

$$\begin{bmatrix} 1.3 \\ 5.1 \\ 2.2 \\ 0.7 \\ 1.1 \end{bmatrix}$$

Softmax
activation function

$$\frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$



Softmax

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

Output
layer

$$\begin{bmatrix} 1.3 \\ 5.1 \\ 2.2 \\ 0.7 \\ 1.1 \end{bmatrix}$$



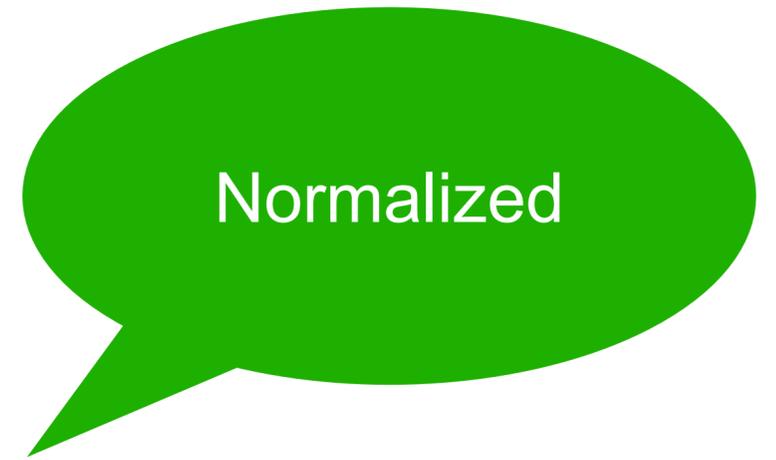
Softmax
activation function

$$\frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$



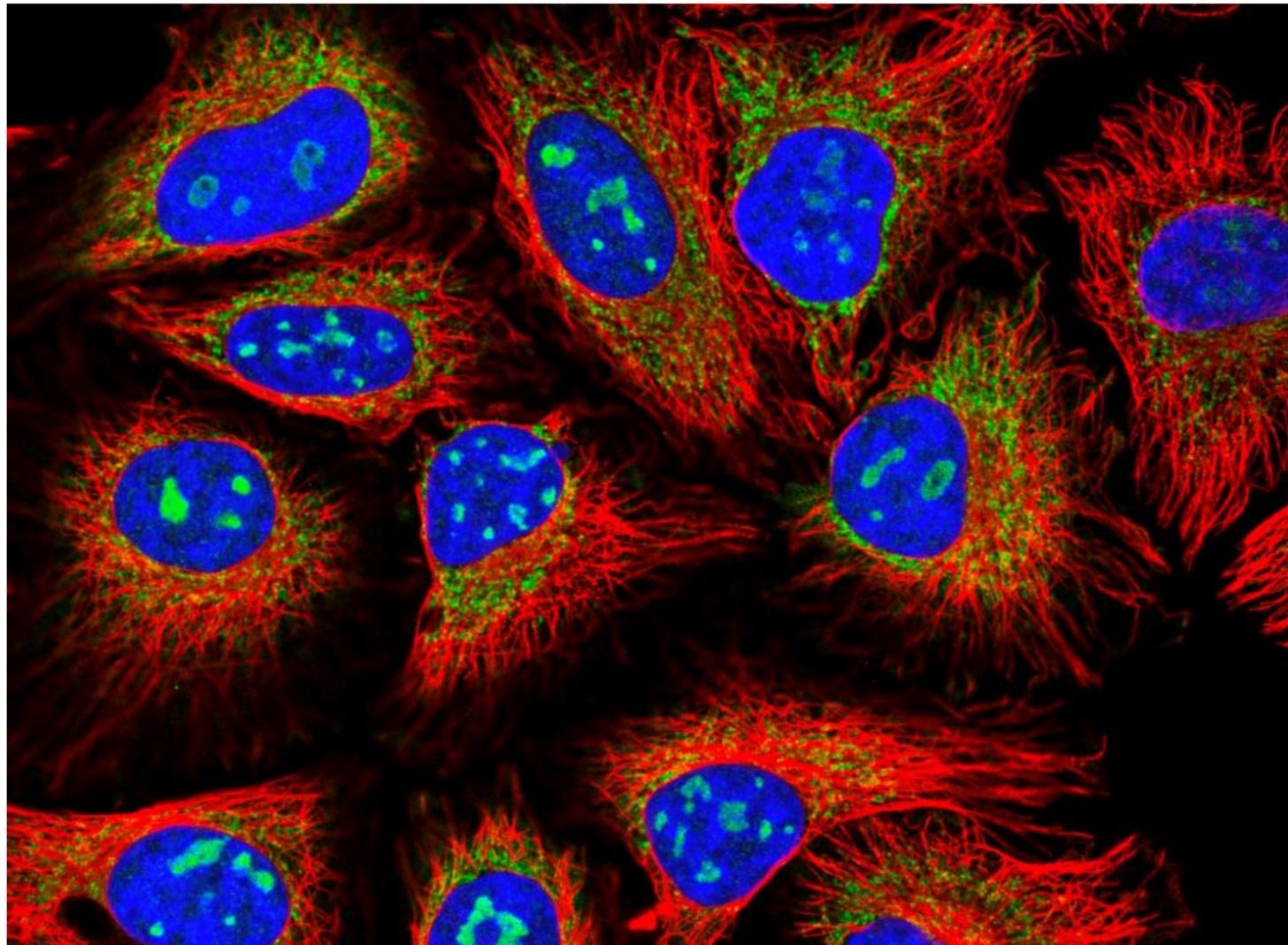
Probabilities

$$\begin{bmatrix} 0.02 \\ 0.90 \\ 0.05 \\ 0.01 \\ 0.02 \end{bmatrix}$$



Classification Tasks at Kaggle

Classify human protein microscope images into 28 categories

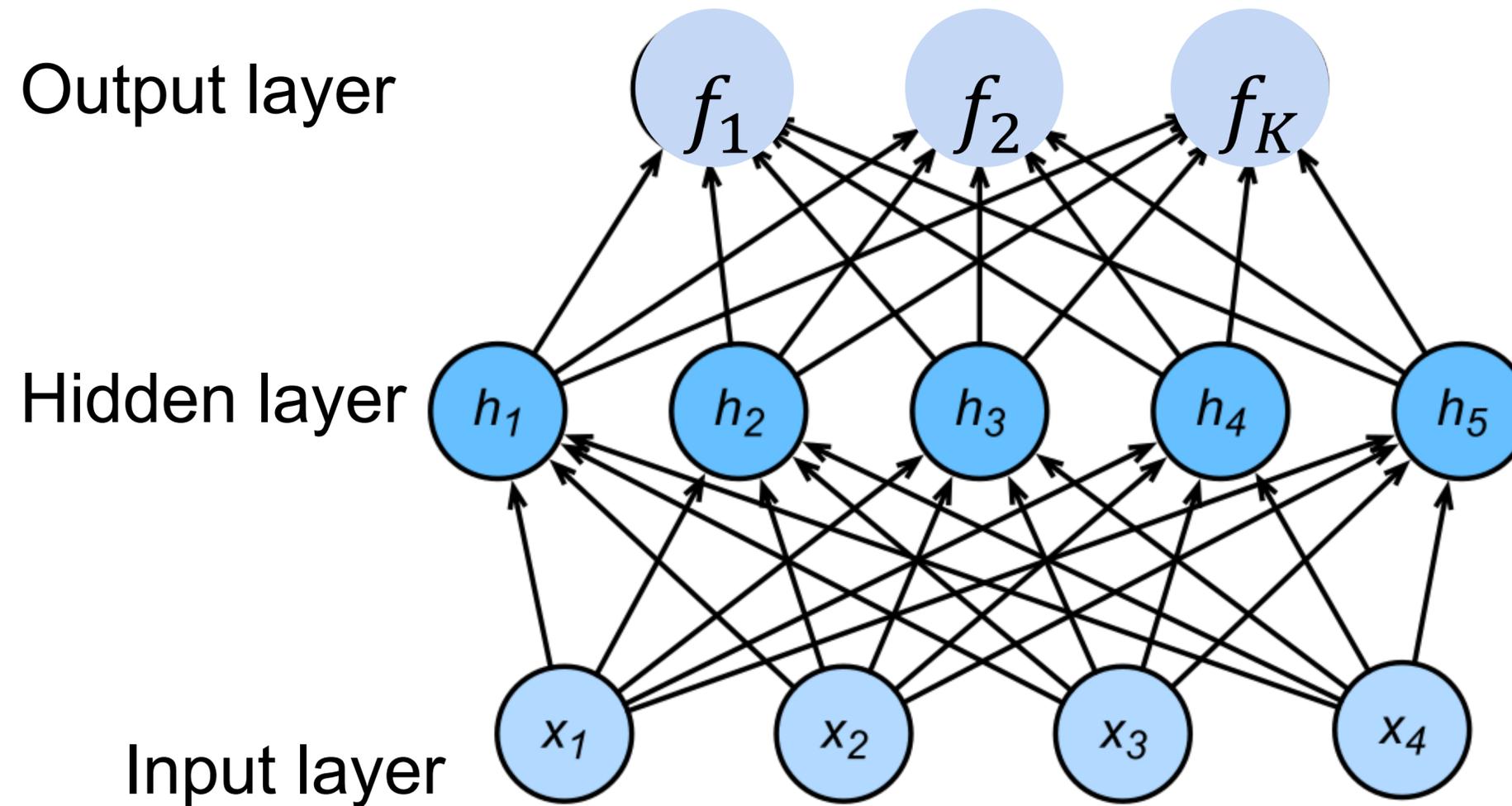


0. Nucleoplasm
1. Nuclear membrane
2. Nucleoli
3. Nucleoli fibrillar
4. Nuclear speckles
5. Nuclear bodies
6. Endoplasmic reticu
7. Golgi apparatus
8. Peroxisomes
9. Endosomes
10. Lysosomes
11. Intermediate fila
12. Actin filaments
13. Focal adhesion si
14. Microtubules
15. Microtubule ends
16. Cytokinetic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

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

More complicated neural networks

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



More complicated neural networks

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

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

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

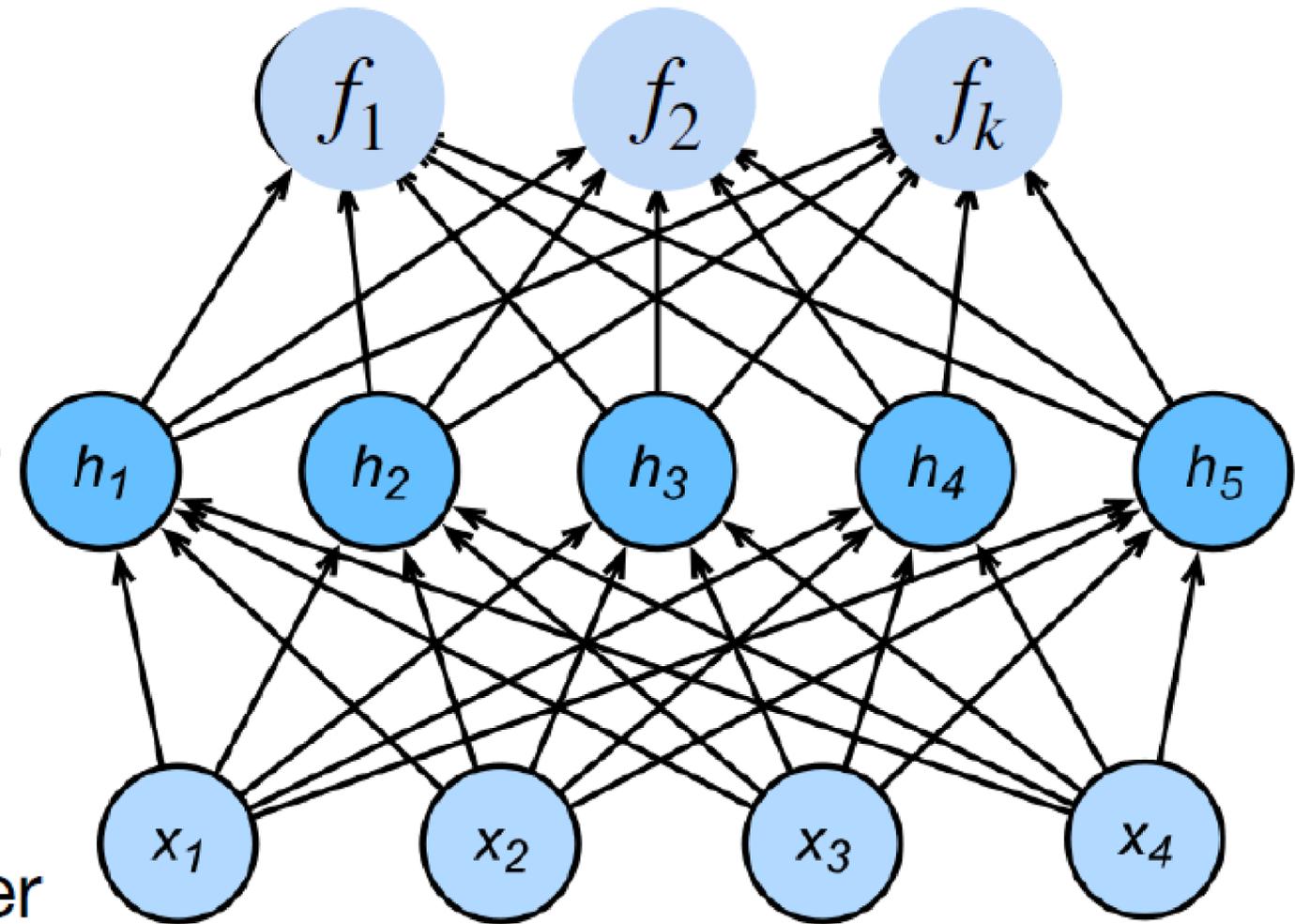
$$\mathbf{f} = \mathbf{W}^{(2)}\mathbf{h} + \mathbf{b}^{(2)}$$

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

Output layer

Hidden layer

Input layer



More complicated neural networks: multiple hidden layers

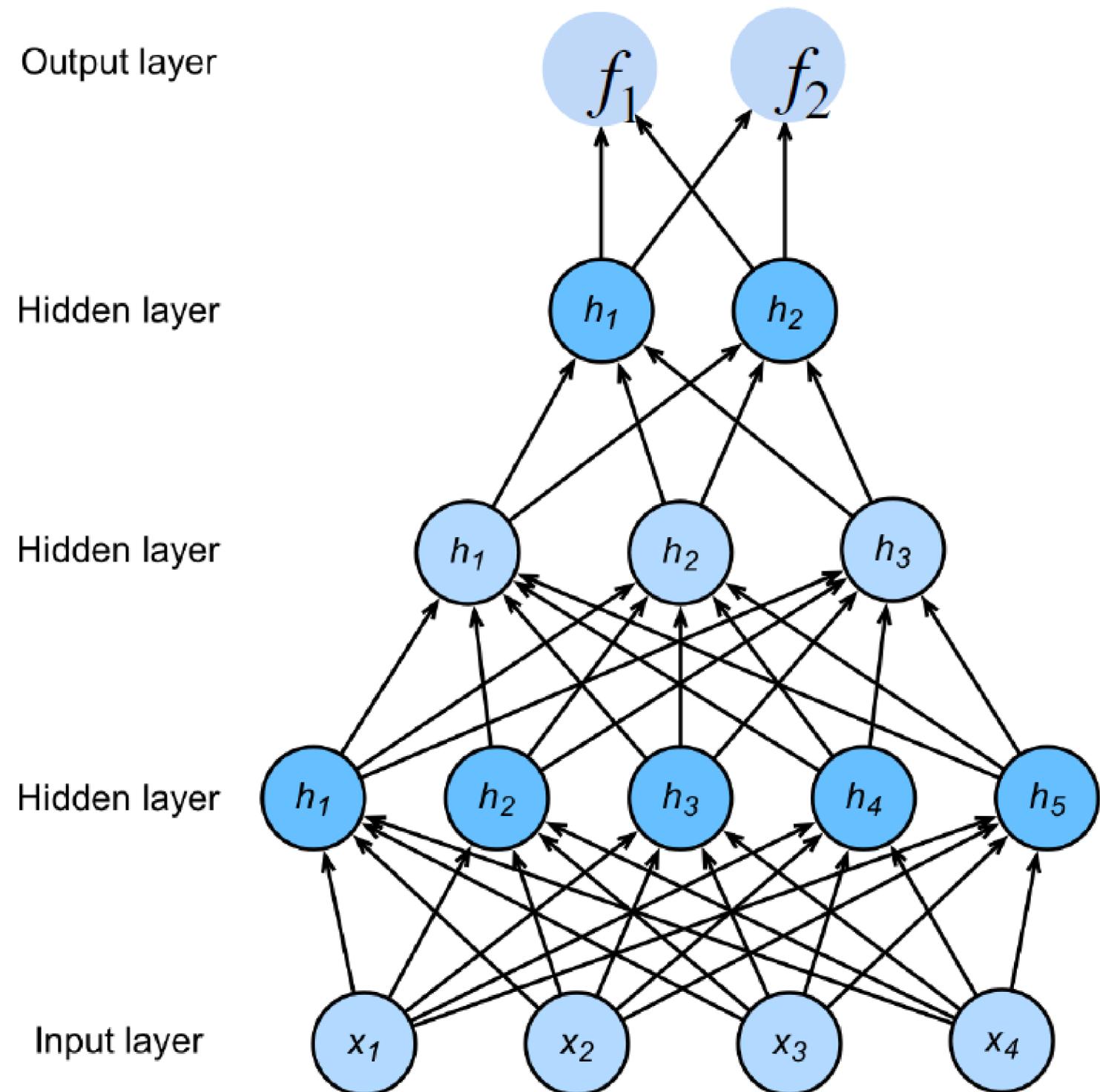
$$\mathbf{h}_1 = \sigma(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})$$

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

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

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

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



Quiz Break

Q1.1

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

Q1.1

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

Q1.2

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

Q1.2

Suppose you are given a 3-layer multilayer perceptron (2 hidden layers h1 and h2 and 1 output layer). All activation functions are sigmoids, and the output layer uses a softmax function. Suppose h1 has 1024 units and h2 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?

$$1024 * 2 + 1024 + 512 * 1024 + 512 + 512 * 3 + 3 = 529411$$

Quiz Break

Q1.3

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

Quiz Break

Q1.3

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

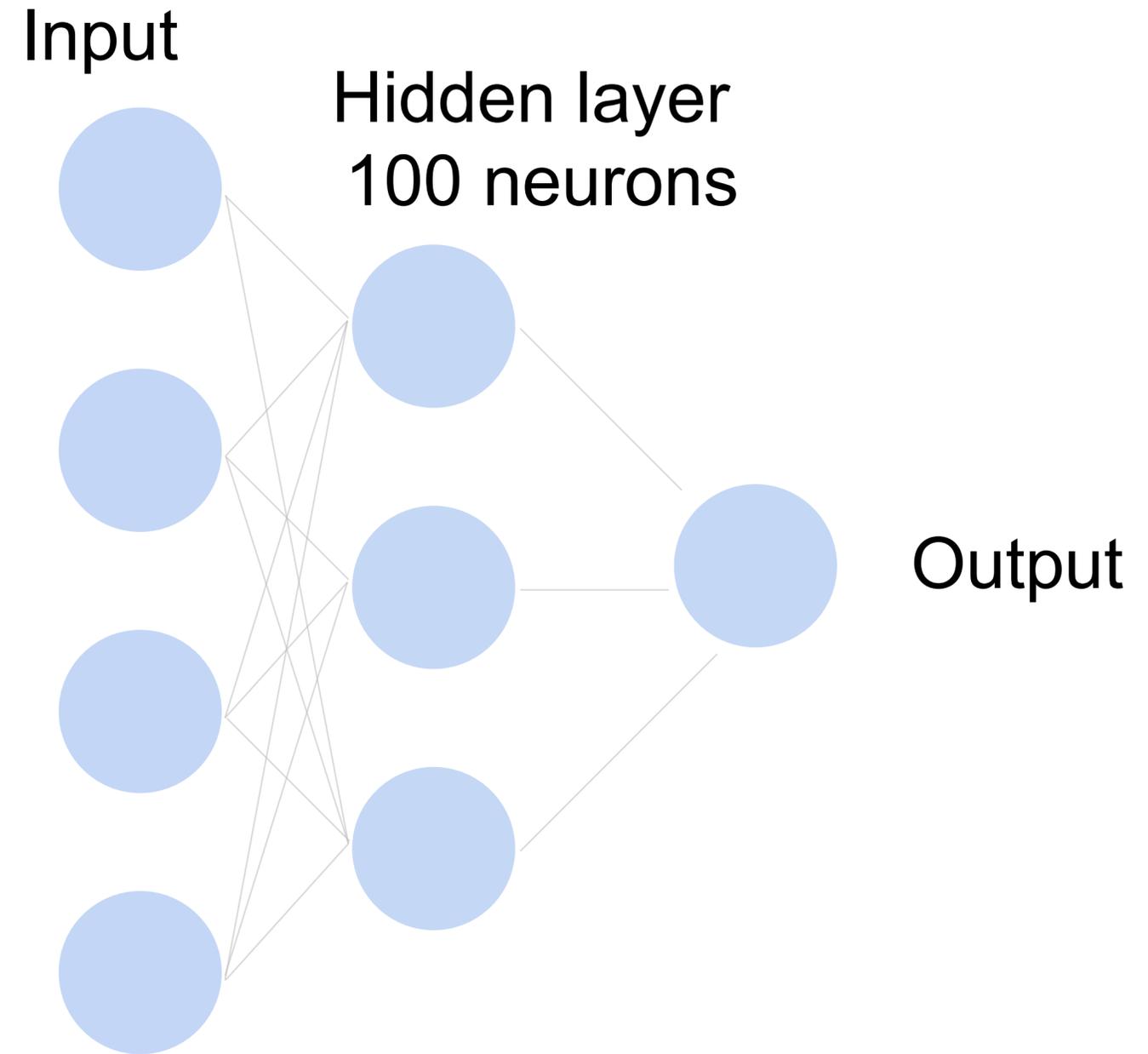


Solution:

A combination of linear Perceptrons is still a linear function.

How to train a neural network?

Classify cats vs. dogs



How to train a neural network? Binary classification

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

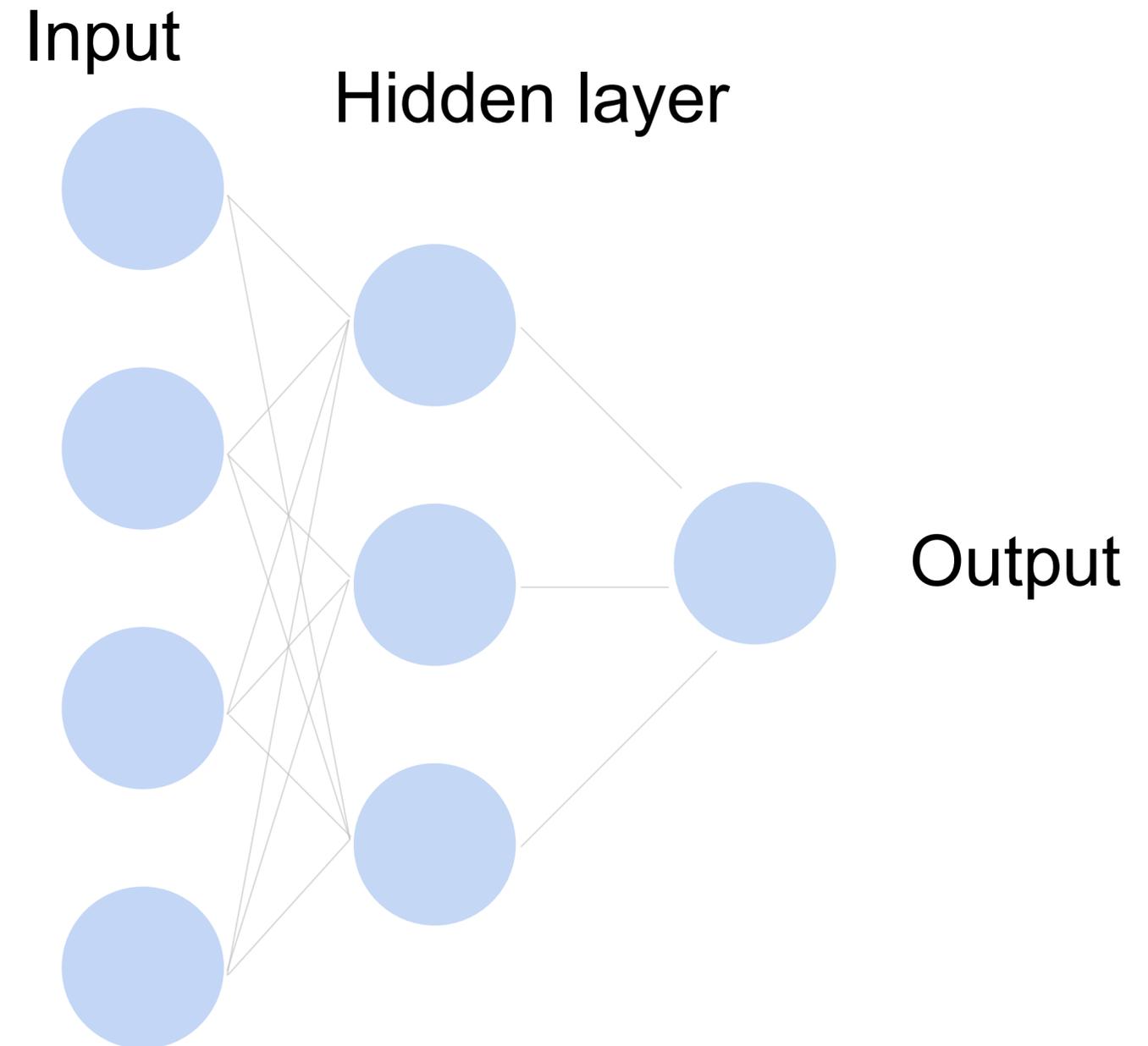
$\hat{y} \in [0,1]$ Model output for example \mathbf{x}

(This is a function of all weights $W: \hat{y} = g(W)$)

y Ground truth label for example \mathbf{x}

Learning by matching the output to the label

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



How to train a neural network? Binary classification

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

Per-sample loss:

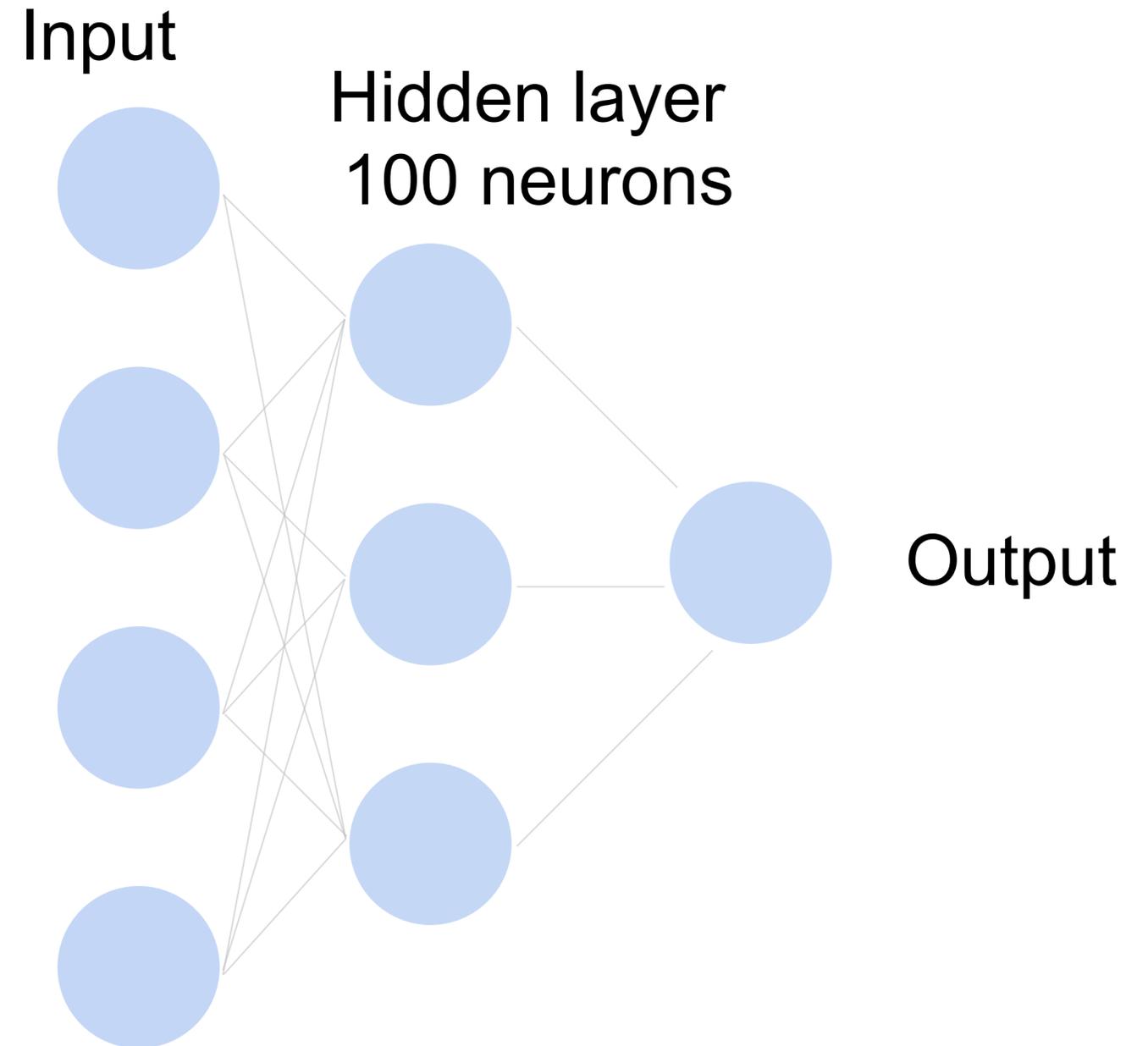
$$\ell(\mathbf{x}, y) = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y})$$



Negative log likelihood

Minimizing NLL is equivalent to Max Likelihood Learning (MLE)

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

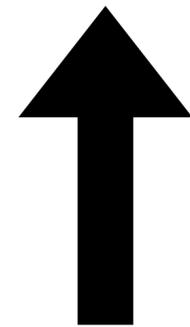


How to train a neural network?

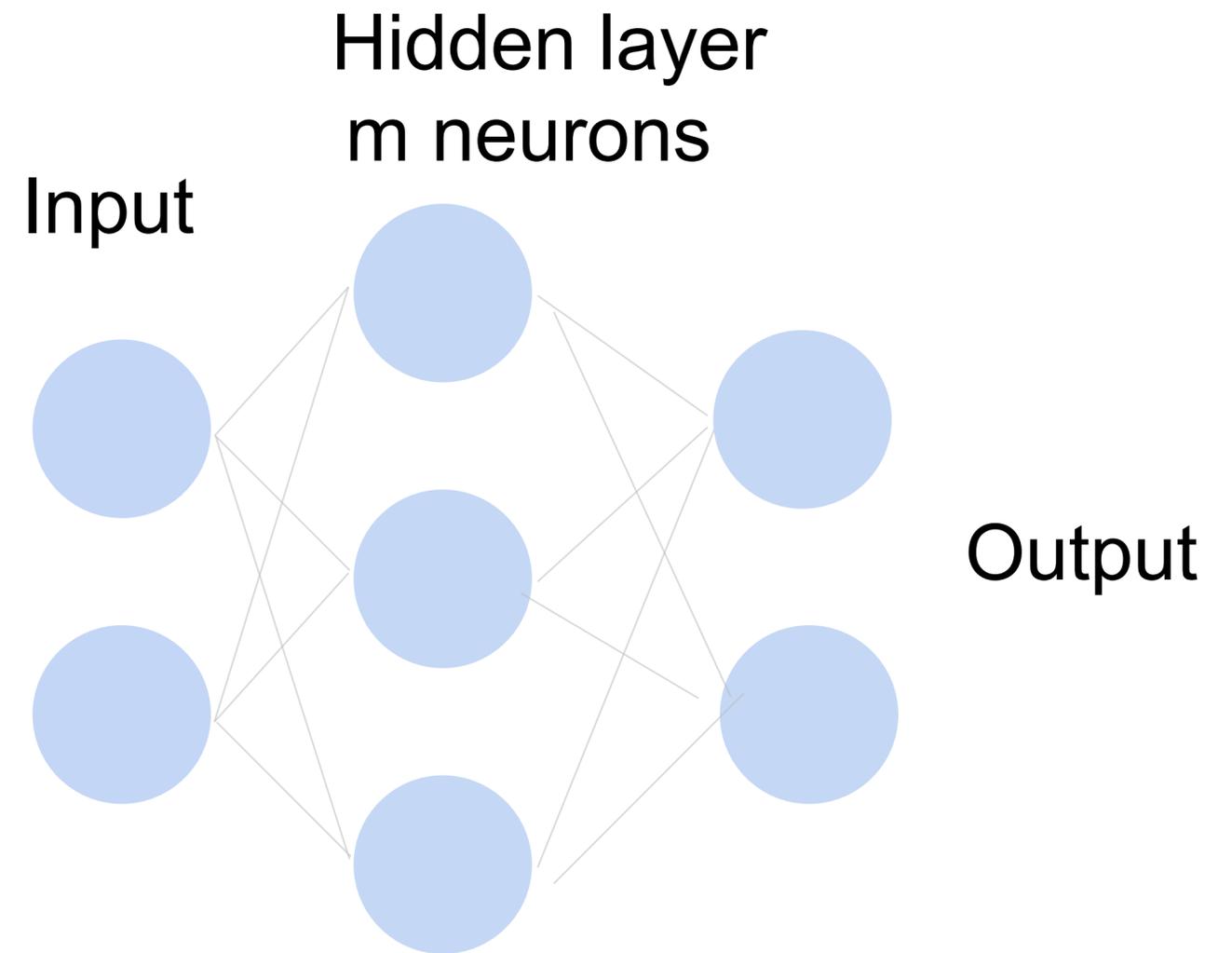
Loss function: $\frac{1}{|D|} \sum_i \ell(\mathbf{x}_i, y_i)$

Per-sample loss:

$$\ell(\mathbf{x}, y) = \sum_{j=1}^K -y_j \log p_j$$



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



Cross-Entropy Loss

softmax
(model prediction)

True label

Neural Networks

0.8

0.2

1

$$L_{CE} = \sum_j -y_j \log(p_j)$$

$$= -\log(0.8)$$

p

Y

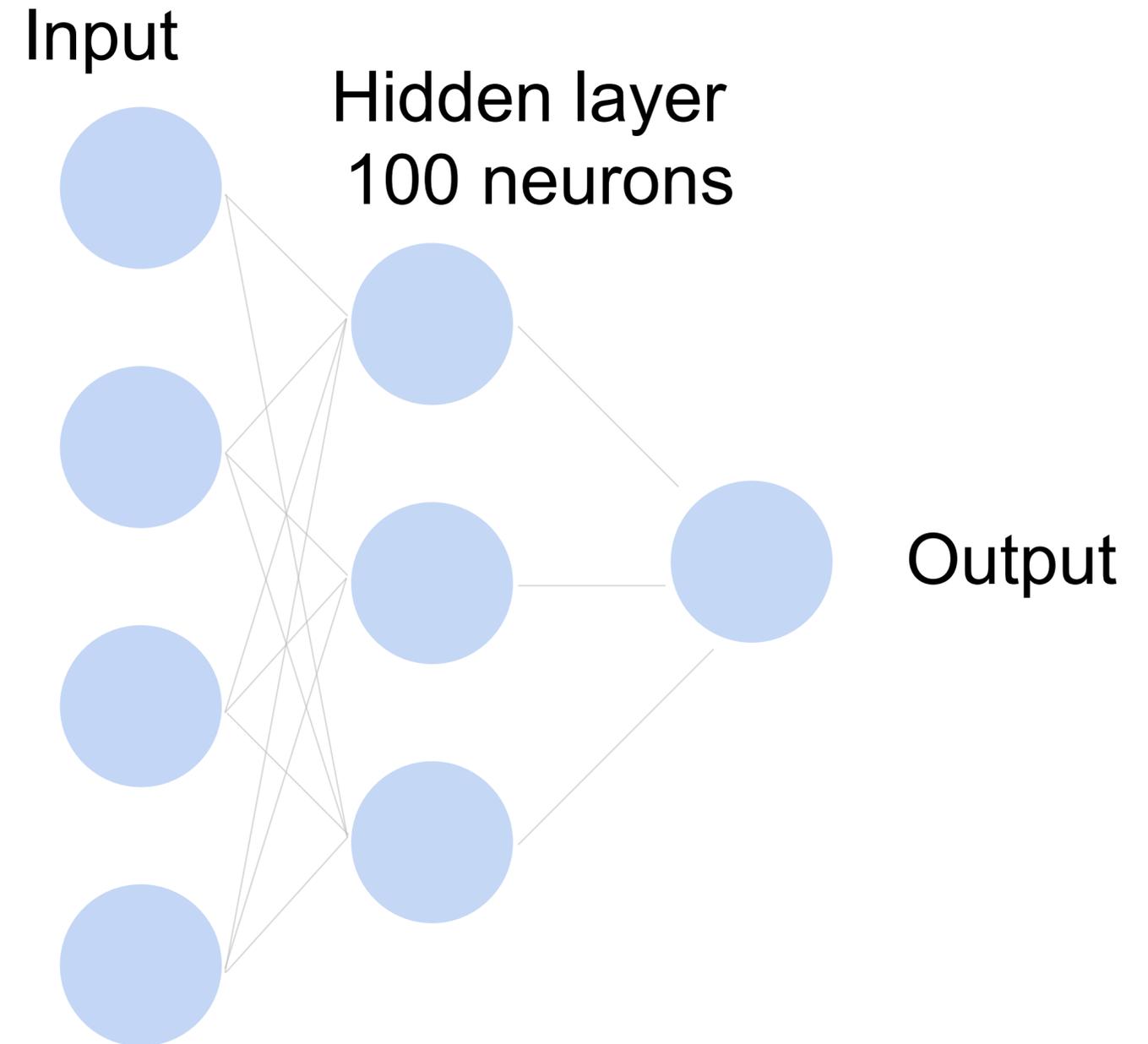
Goal: push p and Y to be identical

How to train a neural network?

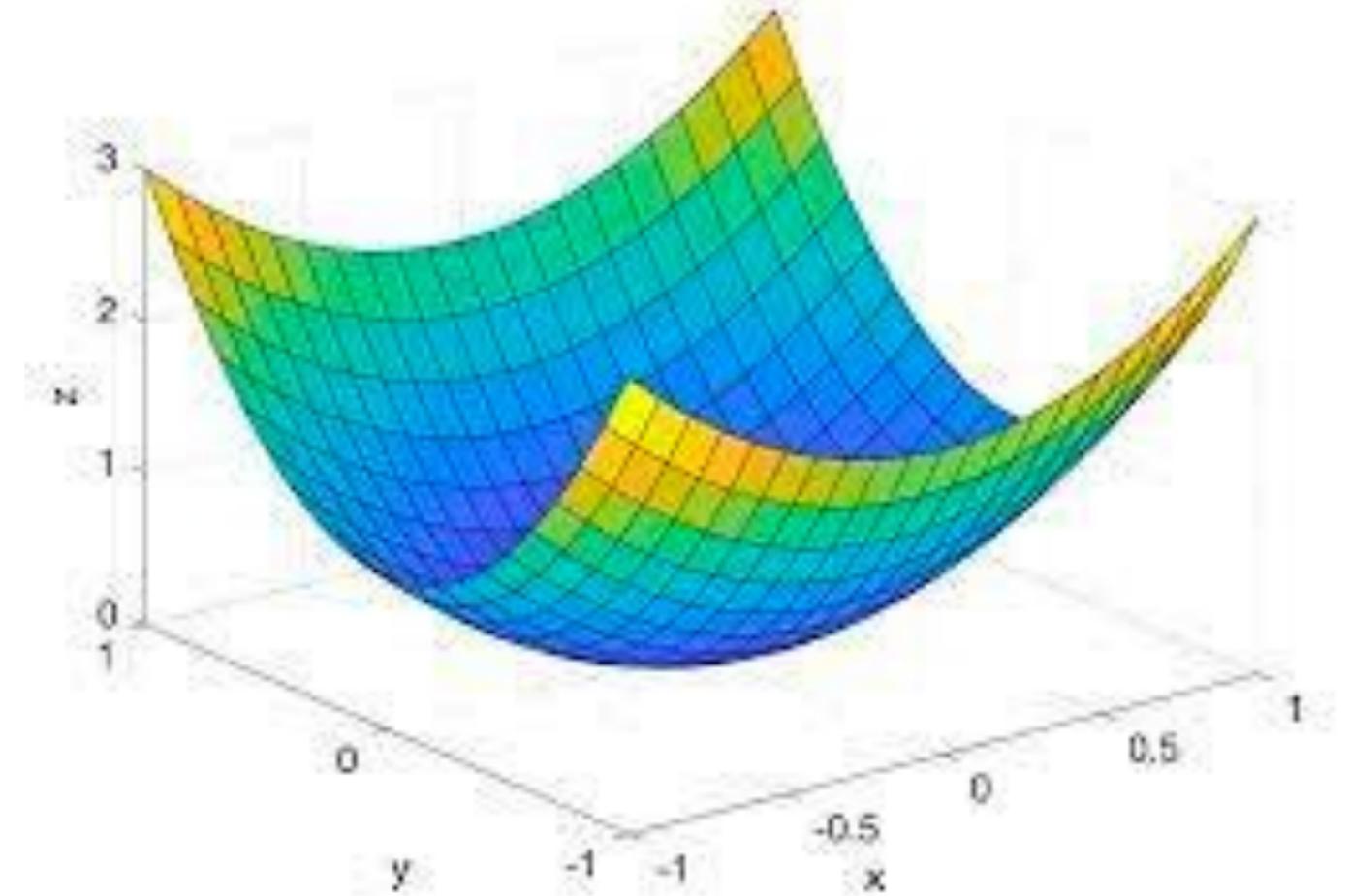
Update the weights W to minimize the loss function

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

Use gradient descent!



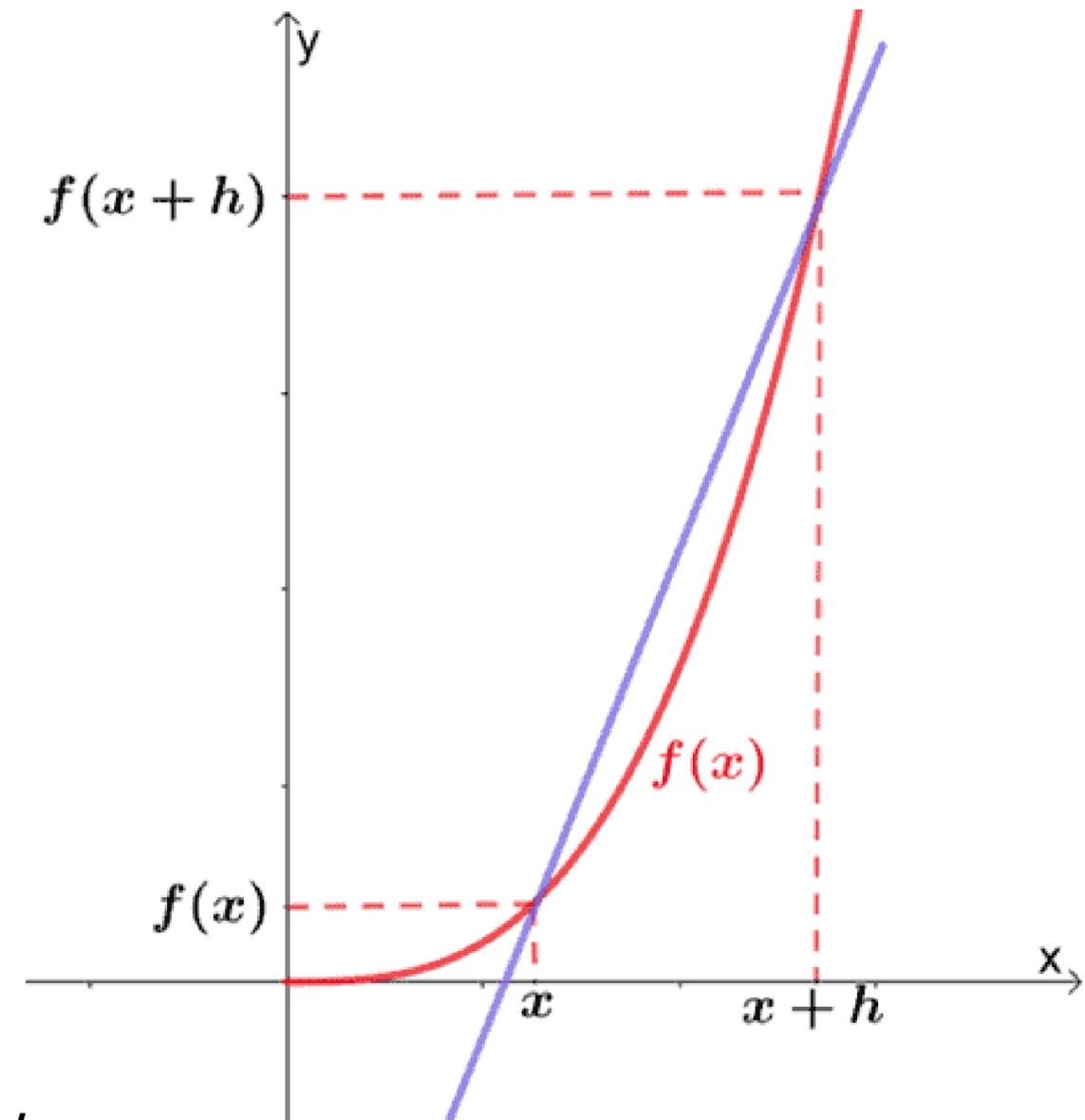
Detour: Multivariate Calculus



Derivatives of functions of single variables

- Given function $f(x)$, we would like to find the slope of the tangent line at any point x_0 .
- Why? Tells us how fast $f(x)$ is increasing / decreasing at x_0 .

- $$\frac{df}{dx} = \lim_{\Delta x \rightarrow 0} \frac{f(x + \Delta x) - f(x)}{\Delta x}$$



Derivatives of functions of single variables

- For many functions $f(x)$ we have simple rules to find the derivative.
- If $f(x) = cx^2$ then $\frac{df}{dx} = 2cx$
 - Or more generally, if $f(x) = cx^p$ then $\frac{df}{dx} = cp x^{p-1}$
- If $f(x) = \log x$ then $\frac{df}{dx} = \frac{1}{x}$
- If $f(x) = c$ then $\frac{df}{dx} = 0$

More Complex Functions

- Derivation rules can be applied hierarchically to find derivatives of more complex functions.
- **Sum rule:** If $f(x) = h(x) + g(x)$ then $\frac{df}{dx} = \frac{dh}{dx} + \frac{dg}{dx}$
- **Product rule:** If $f(x) = h(x) \cdot g(x)$ then $\frac{df}{dx} = h(x) \frac{dg}{dx} + g(x) \frac{dh}{dx}$
- **Chain rule:** If $f(x) = h(g(x))$ then $\frac{df}{dx} = \frac{dh}{dg} \frac{dg}{dx}$
- **More complex chain rule:** if $f(x) = f_1(f_2(\cdots f_n(x) \cdots))$ then $\frac{df}{dx} = \frac{df_1}{df_2} \frac{df_2}{df_3} \cdots \frac{df_n}{dx}$

Derivatives of functions of multiple variables

- Generalize derivative to functions of form $f(x_1, x_2, \dots, x_n)$.
- The partial derivative $\frac{\partial f}{\partial x_1}$ tells us how fast f increases / decreases as we increase the value of x_1 .
- For each variable x_i we have a partial derivative $\frac{\partial f}{\partial x_i}$.
- The **gradient** is the vector of all of these partial derivatives.
 - $\frac{df}{d\mathbf{x}} = \left[\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right]$

Computing Partial Derivatives

- Rules from single-variable calculus directly extend to multivariate calculus with one small change.
 - When computing $\frac{\partial f}{\partial x_i}$, treat all other variables as constants.
- Examples:
 - If $f(x_1, x_2) = \log x_1 + \log x_2$ then $\frac{\partial f}{\partial x_1} = \frac{1}{x_1}$
 - If $f(x_1, x_2) = x_1(x_1 + x_2)$ then $\frac{\partial f}{\partial x_2} = x_1$

Quiz Break

Q2.1 What is the derivative of $f(x) = \log(1 - x)$

A. $\log(1 - x)$

B. $-\frac{1}{x}$

C. $\frac{1}{(1-x)}$

D. $-\frac{1}{(1-x)}$

Quiz Break

Q2.1 What is the derivative of $f(x) = \log(1 - x)$

A. $\log(1 - x)$

B. $-\frac{1}{x}$

C. $\frac{1}{(1-x)}$

D. $-\frac{1}{(1-x)}$

Let $a = 1 - x$

$$\frac{df}{dx} = \frac{df}{da} \frac{da}{dx}$$

$$\frac{df}{da} = \frac{1}{a} \quad \frac{da}{dx} = -1$$

$$\frac{df}{dx} = \frac{1}{a} (-1) = -\frac{1}{(1-x)}$$

Quiz Break

Q2.2 What is the derivative of $\frac{\partial f}{\partial w_1}$ of: $f(x_1, x_2, w_1, w_2) = w_1 x_1 + w_2 x_2$

- A. w_1
- B. x_1
- C. $x_1 + w_2 x_2$
- D. $w_1 + w_2 x_2$

Quiz Break

Q2.2 What is the derivative of $\frac{\partial f}{\partial w_1}$ of: $f(x_1, x_2, w_1, w_2) = w_1x_1 + w_2x_2$

- A. w_1
- B. x_1**
- C. $x_1 + w_2x_2$
- D. $w_1 + w_2x_2$

We treat all other variables as constants
and we take partial derivative w.r.t. w_1

Quiz Break

Q2.3 What is the partial derivative $\frac{\partial f}{\partial w_1}$ of:

$$f(x_1, x_2, w_1, w_2, y) = y \log \sigma(w_1 x_1 + w_2 x_2) + (1 - y) \log(1 - \sigma(w_1 x_1 + w_2 x_2))$$

when $y = 1$ and $\sigma(z) = \frac{1}{1+e^{-z}}$. **Hint:** $\frac{\partial \sigma}{\partial z} = \sigma(z)(1 - \sigma(z))$.

- A. $(1 - \sigma(w_1 x_1 + w_2 x_2))x_1$
- B. $(1 - \sigma(w_1 x_1 + w_2 x_2))$
- C. $\sigma(w_1 x_1 + w_2 x_2)x_1$
- D. $\sigma(w_1 x_1 + w_2 x_2)(1 - \sigma(w_1 x_1 + w_2 x_2))x_1$

Quiz Break

Q2.3 What is the partial derivative $\frac{\partial f}{\partial w_1}$ of:

$$f(x_1, x_2, w_1, w_2, y) = y \log \sigma(w_1 x_1 + w_2 x_2) + (1 - y) \log(1 - \sigma(w_1 x_1 + w_2 x_2))$$

when $y = 1$ and $\sigma(z) = \frac{1}{1+e^{-z}}$. **Hint:** $\frac{\partial \sigma}{\partial z} = \sigma(z)(1 - \sigma(z))$.

- A. $(1 - \sigma(w_1 x_1 + w_2 x_2))x_1$
- B. $(1 - \sigma(w_1 x_1 + w_2 x_2))$
- C. $\sigma(w_1 x_1 + w_2 x_2)x_1$
- D. $\sigma(w_1 x_1 + w_2 x_2)(1 - \sigma(w_1 x_1 + w_2 x_2))x_1$

$$\text{Let } z = w_1 x_1 + w_2 x_2 \quad \text{Let } a = \sigma(z)$$

$$\frac{\partial f}{\partial w_1} = \frac{\partial f}{\partial a} \frac{\partial a}{\partial z} \frac{\partial z}{\partial w_1} \quad \text{We set } y = 1$$

$$\frac{\partial f}{\partial a} = \frac{1}{a} \quad \frac{\partial a}{\partial z} = \sigma(z)(1 - \sigma(z)) \quad \frac{\partial z}{\partial w_1} = x_1$$

$$\frac{\partial f}{\partial w_1} = \frac{\partial f}{\partial a} \frac{\partial a}{\partial z} \frac{\partial z}{\partial w_1} = \frac{1}{a} \sigma(z)(1 - \sigma(z))x_1 = (1 - \sigma(w_1 x_1 + w_2 x_2))x_1$$

Gradient Descent

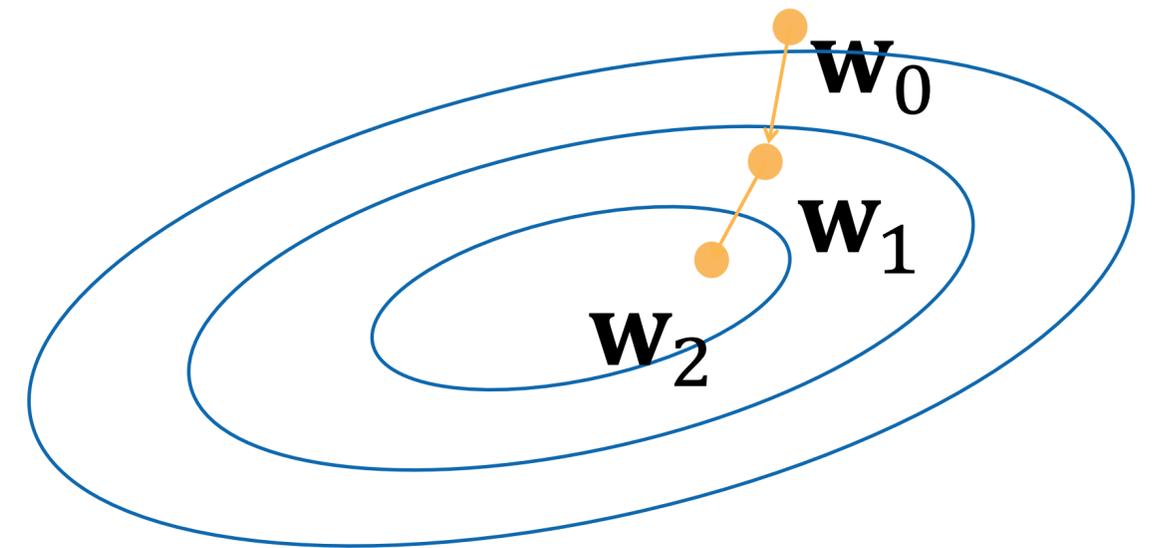
- Choose a learning rate $\eta > 0$
- Initialize the model parameters w_0
- For $t = 1, 2, \dots$
 - Update parameters:

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

D can be very large. Expensive per iteration

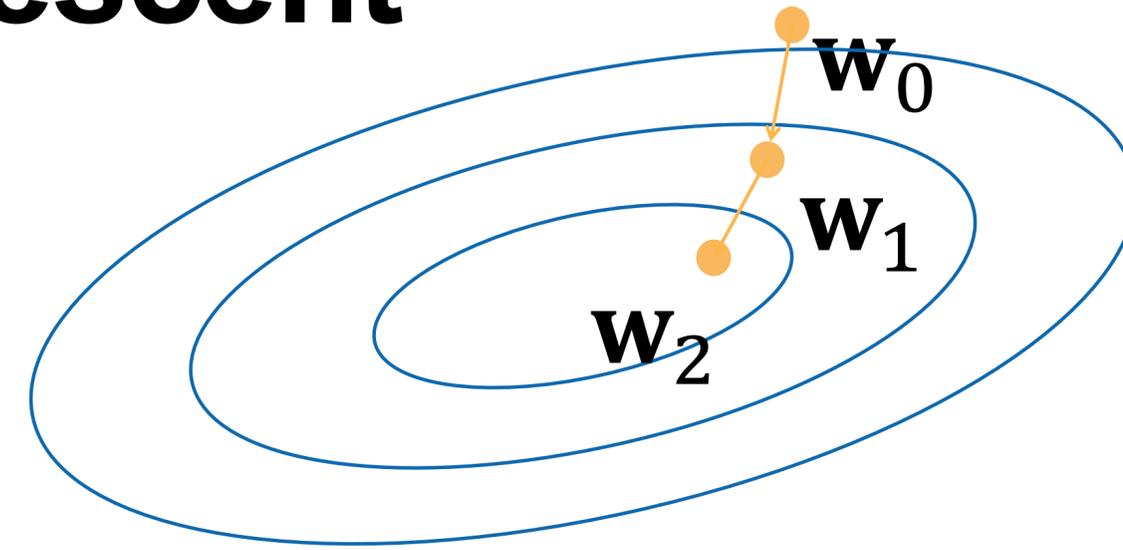
The gradient w.r.t. all parameters is obtained by concatenating the partial derivatives w.r.t. each parameter

- Repeat until converges



Minibatch Stochastic Gradient Descent

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

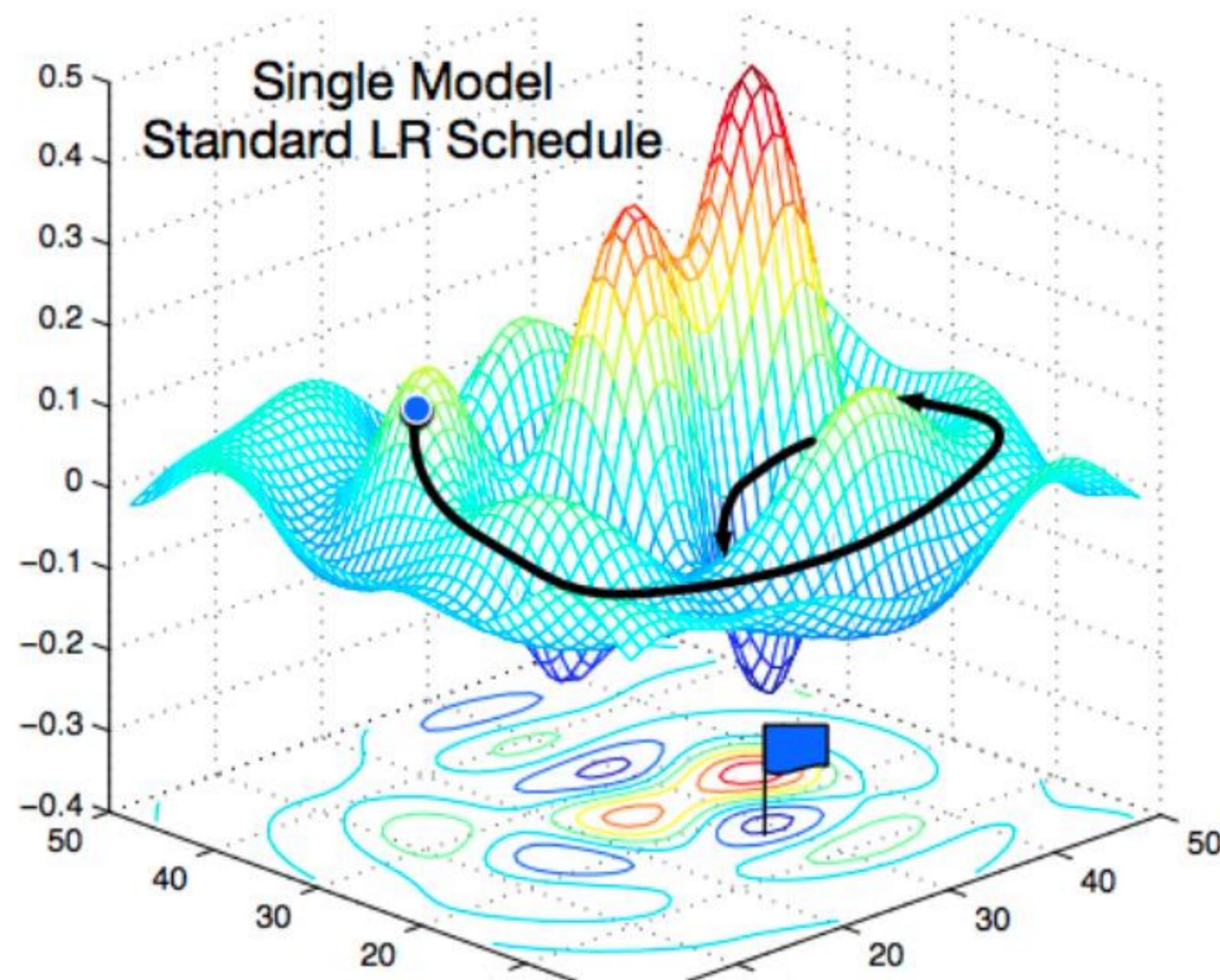


- Randomly sample a subset (mini-batch) $B \subset D$
- Update parameters:

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

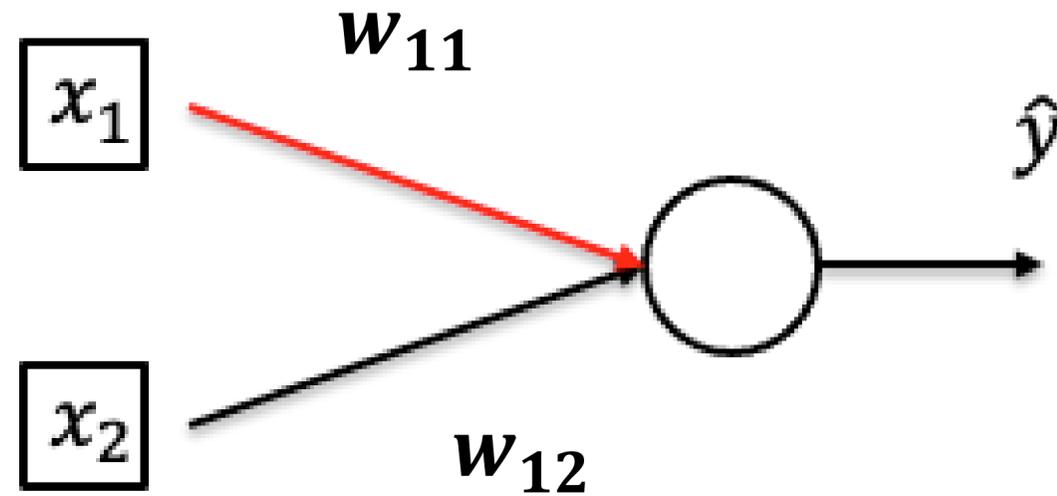
- Repeat until converges

Non-convex Optimization



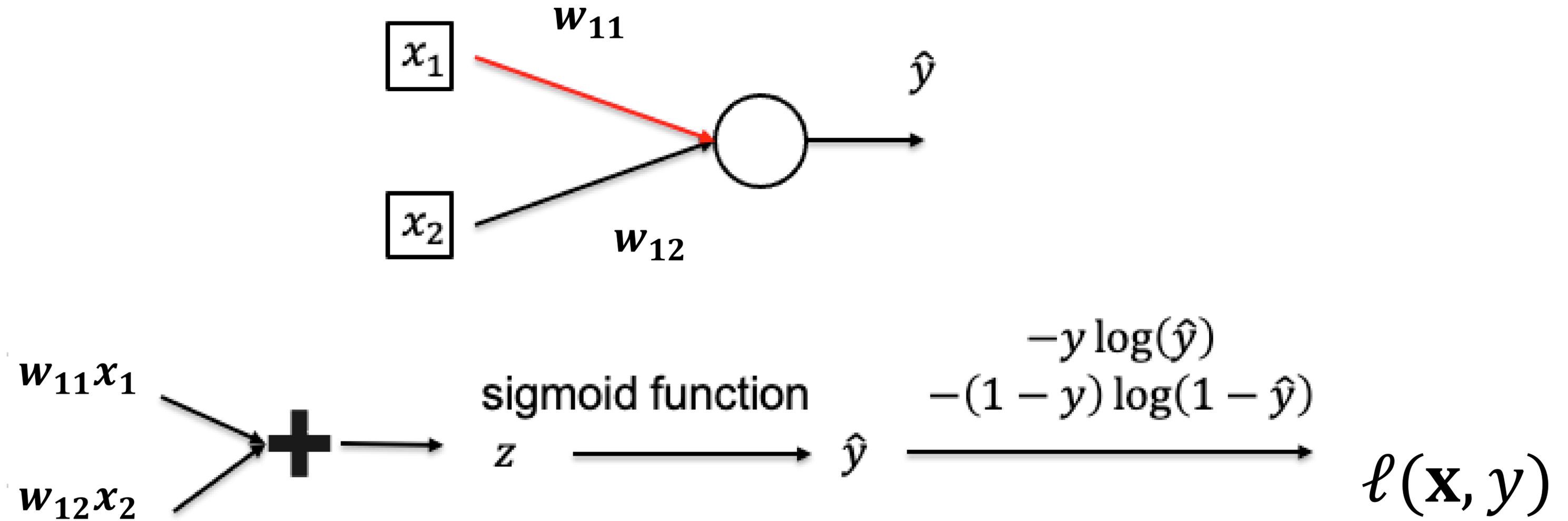
[Gao and Li et al., 2018]

Calculate Gradient (on one data point)



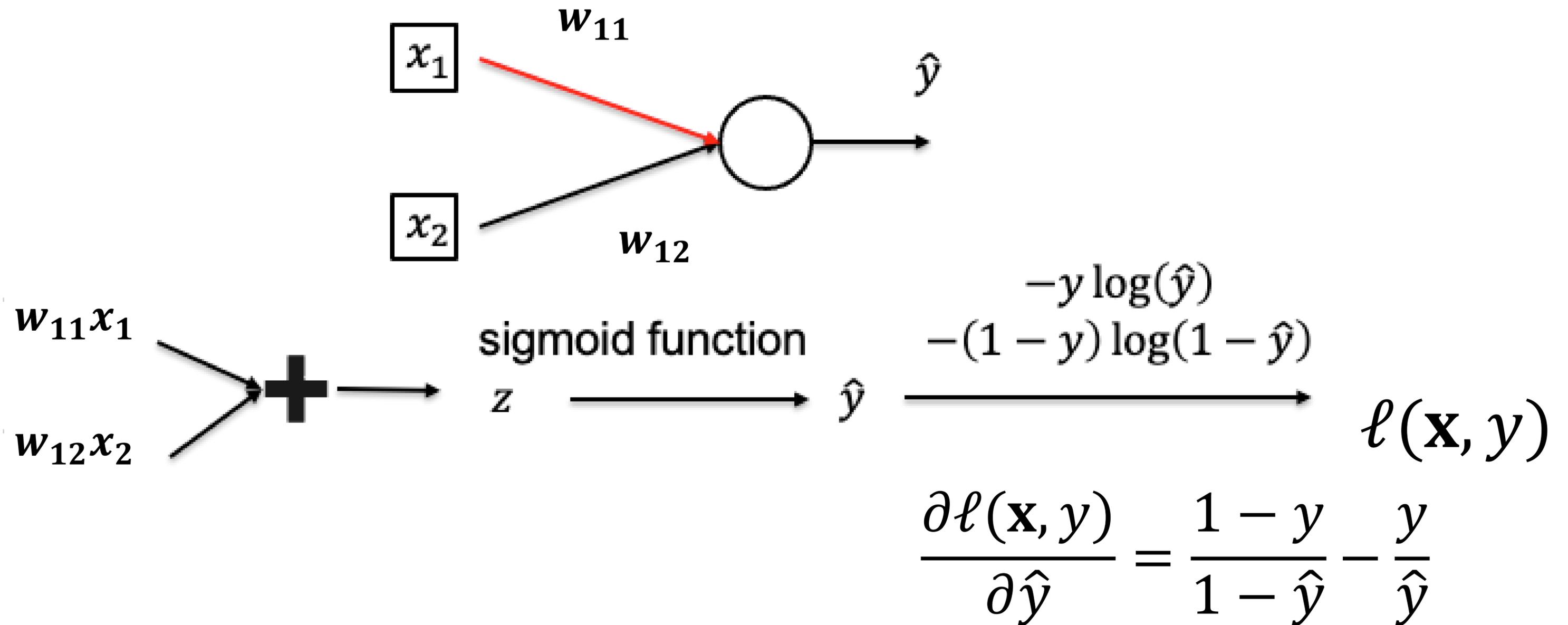
- Want to compute $\frac{\partial \ell(\mathbf{x}, y)}{\partial w_{11}}$
- Data point: $((x_1, x_2), y)$

Calculate Gradient (on one data point)



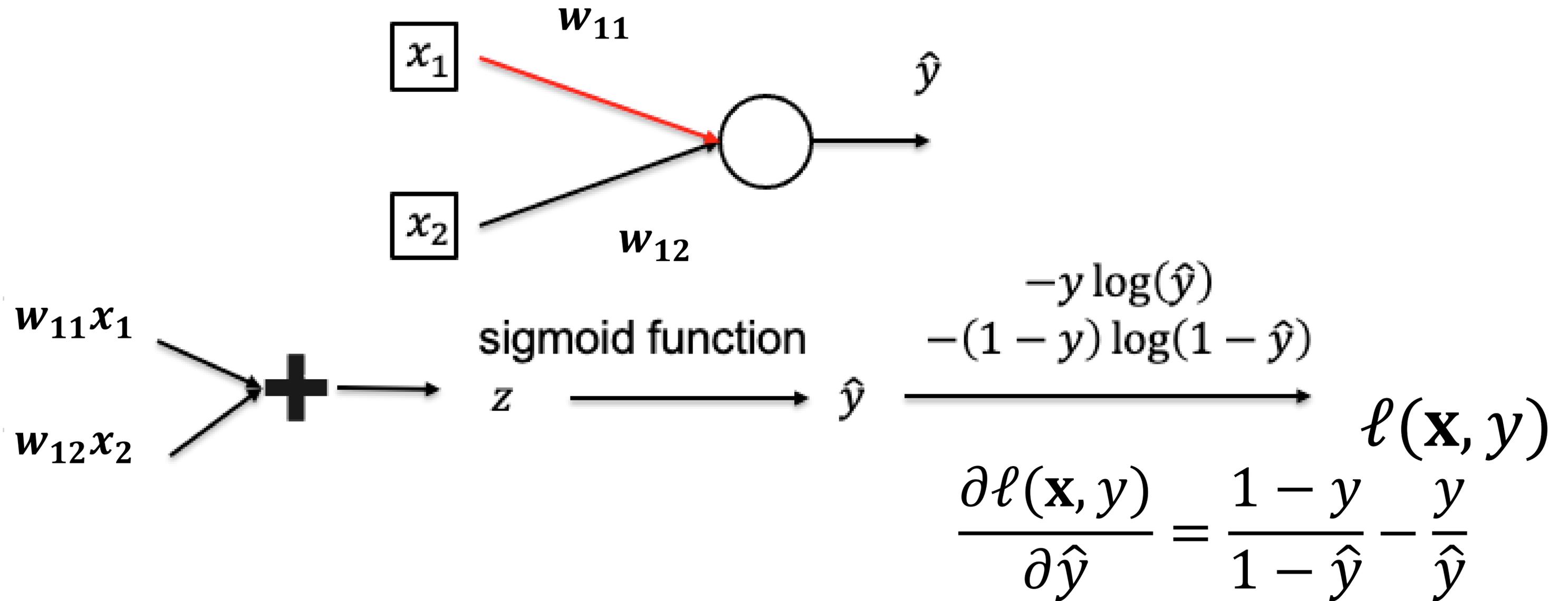
Use chain rule!

Calculate Gradient (on one data point)



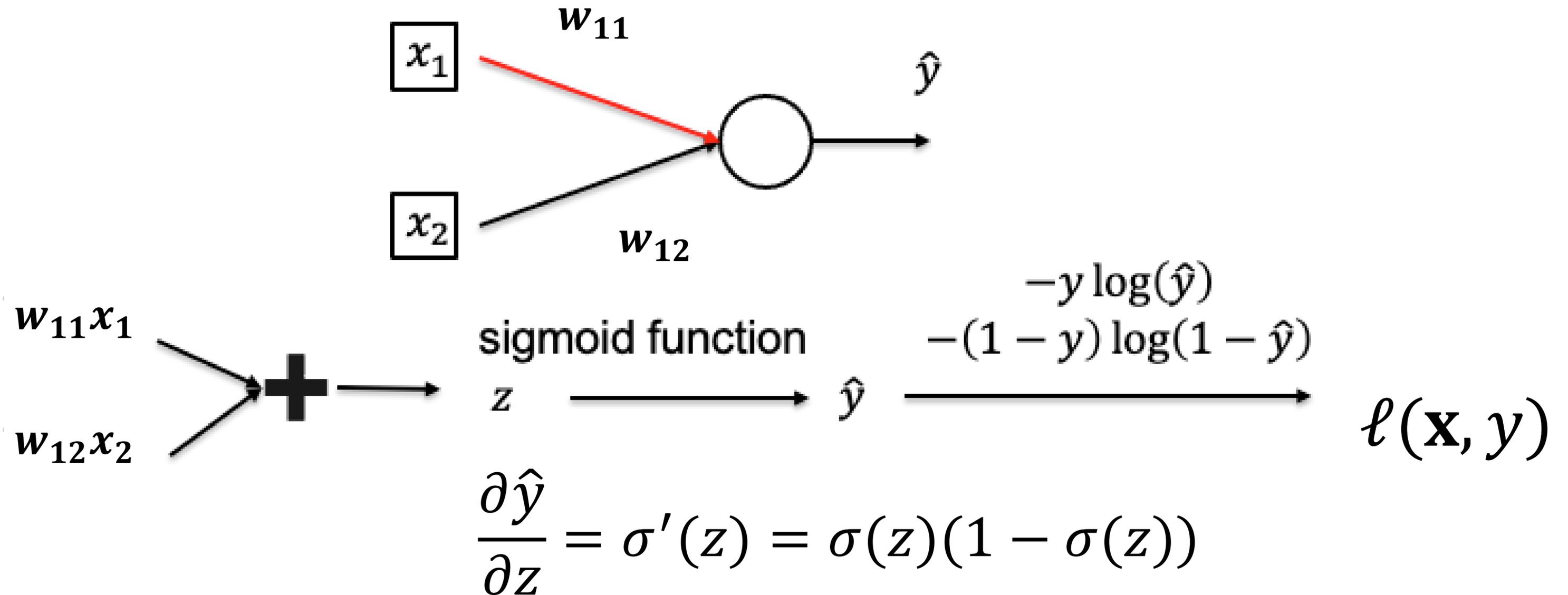
- By chain rule: $\frac{\partial \ell}{\partial w_{11}} = \frac{\partial \ell}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial w_{11}}$

Calculate Gradient (on one data point)



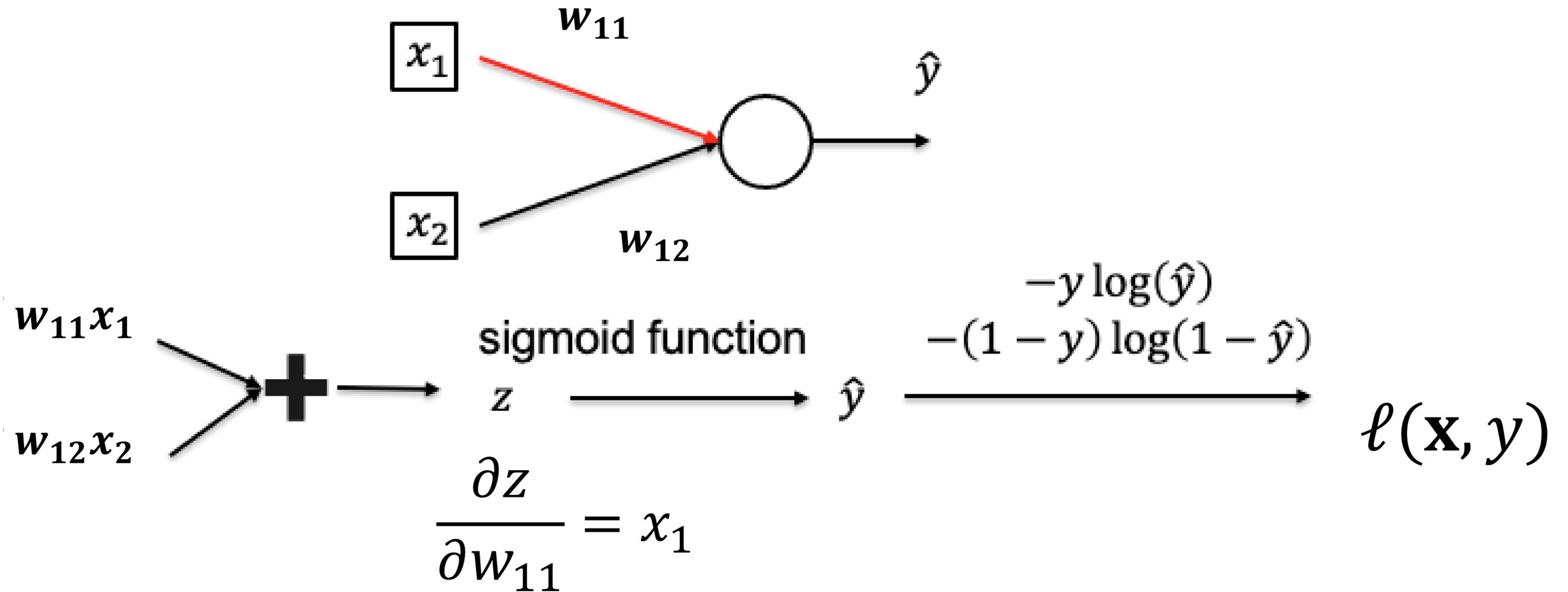
- By chain rule:
$$\frac{\partial \ell}{\partial w_{11}} = \left(\frac{1 - y}{1 - \hat{y}} - \frac{y}{\hat{y}} \right) \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial w_{11}}$$

Calculate Gradient (on one data point)



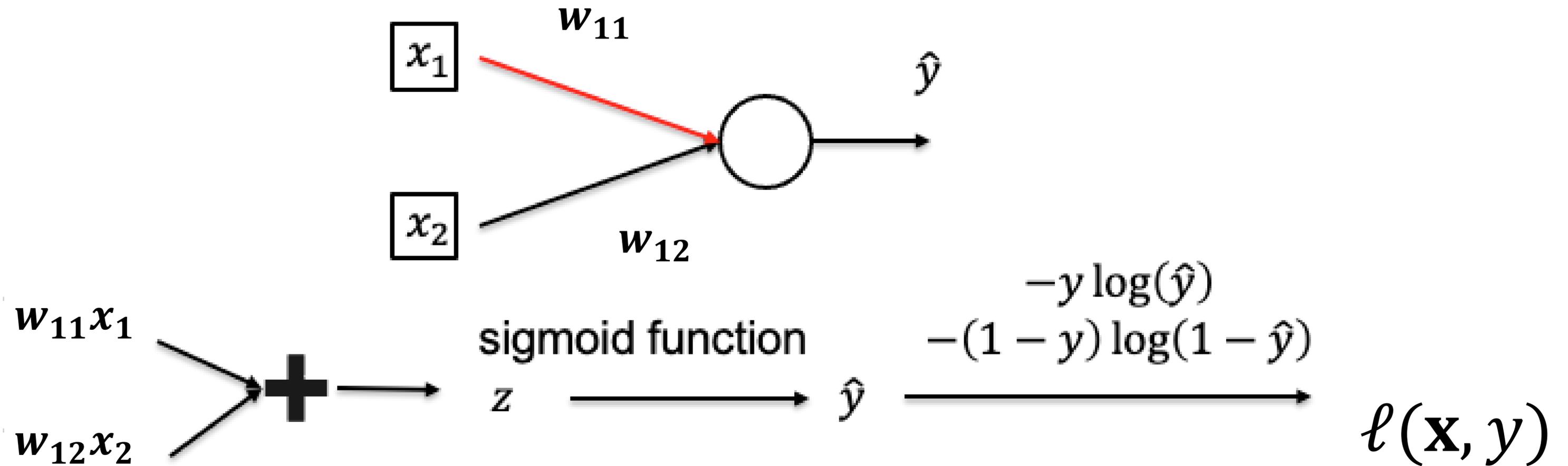
• By chain rule:
$$\frac{\partial l}{\partial w_{11}} = \left(\frac{1-y}{1-\hat{y}} - \frac{y}{\hat{y}} \right) \hat{y} (1-\hat{y}) \frac{\partial z}{\partial w_{11}}$$

Calculate Gradient (on one data point)



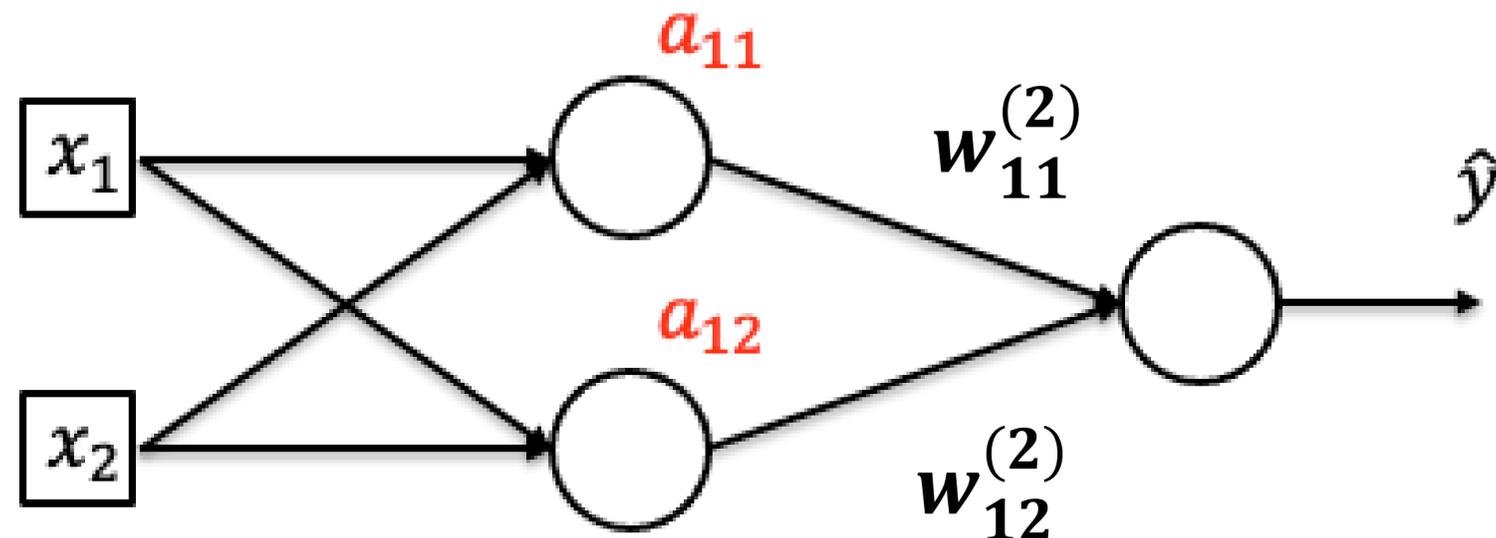
- By chain rule:
$$\frac{\partial \ell}{\partial w_{11}} = \left(\frac{1 - y}{1 - \hat{y}} - \frac{y}{\hat{y}} \right) \hat{y} (1 - \hat{y}) x_1$$

Calculate Gradient (on one data point)

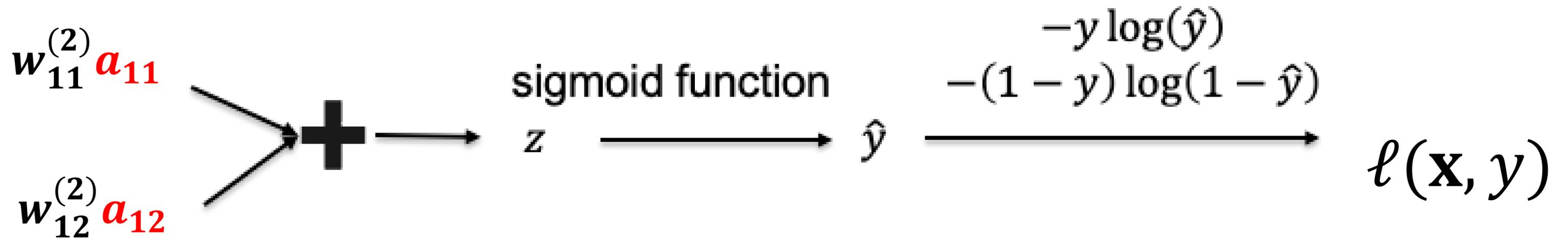


- By chain rule:
$$\frac{\partial \ell}{\partial w_{11}} = (\hat{y} - y)x_1$$

Calculate Gradient (on one data point)

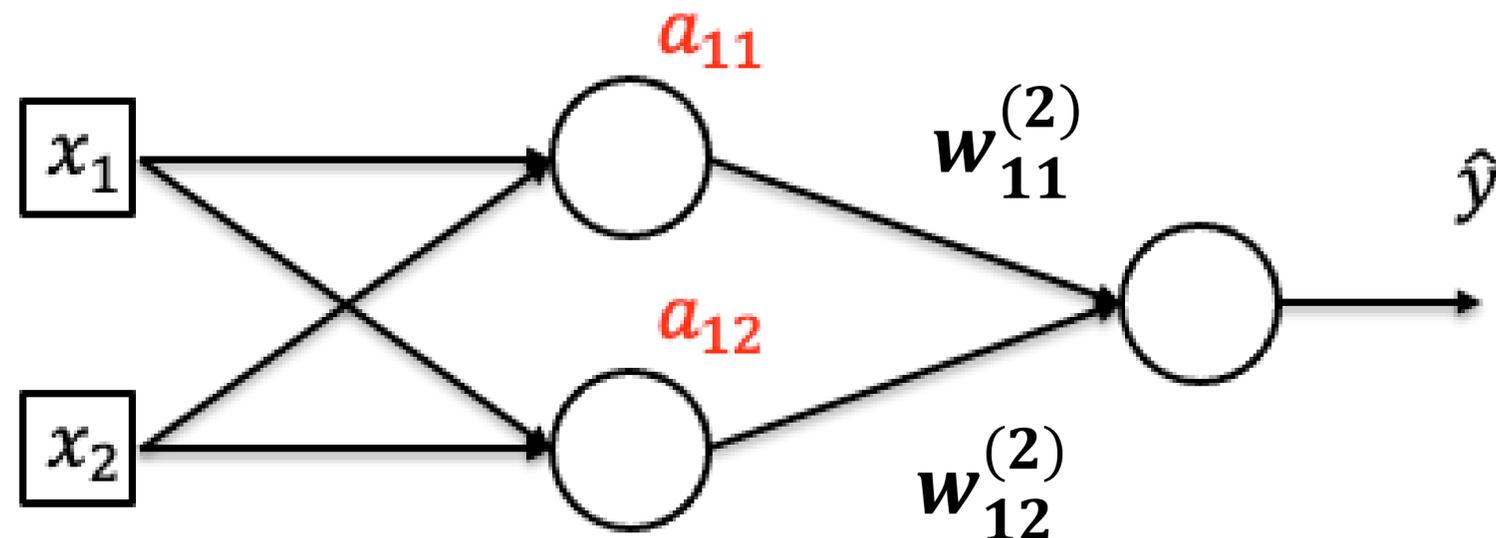


Make it deeper

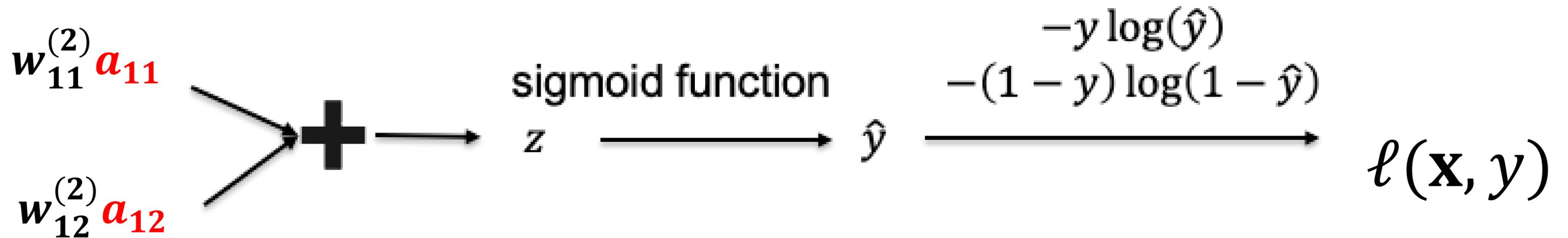


- By chain rule: $\frac{\partial \ell}{\partial a_{11}} = \frac{\partial \ell}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial a_{11}} = (\hat{y} - y) w_{11}^{(2)}$

Calculate Gradient (on one data point)

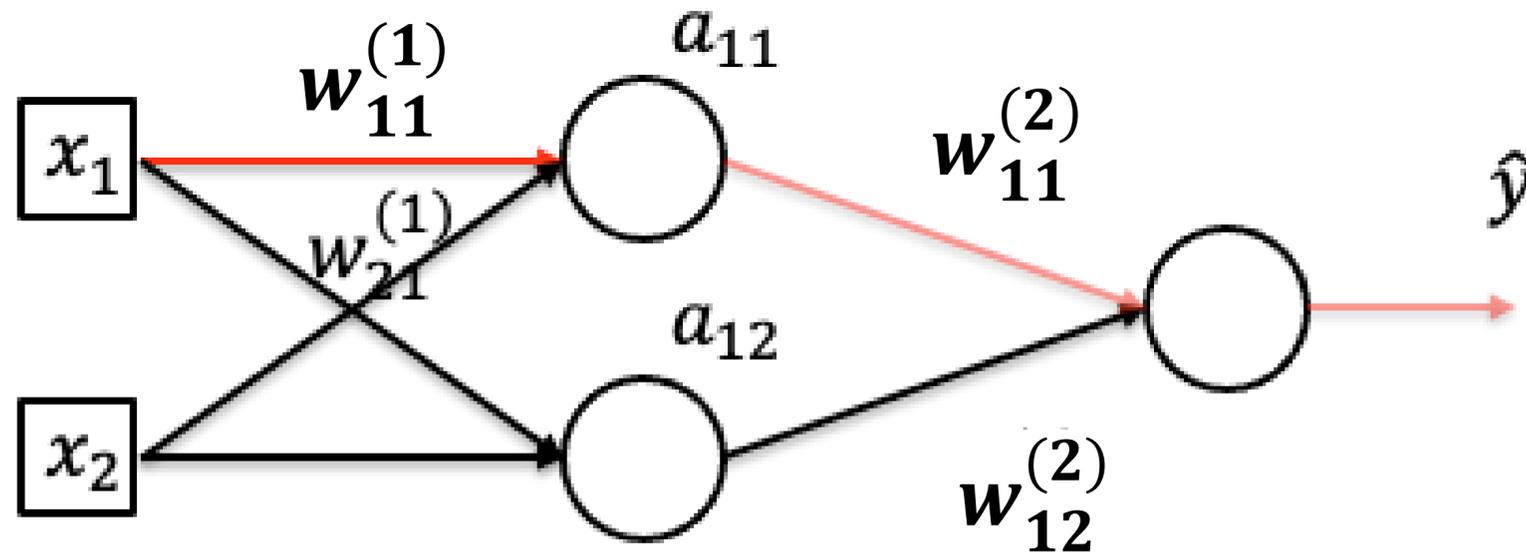


Make it deeper



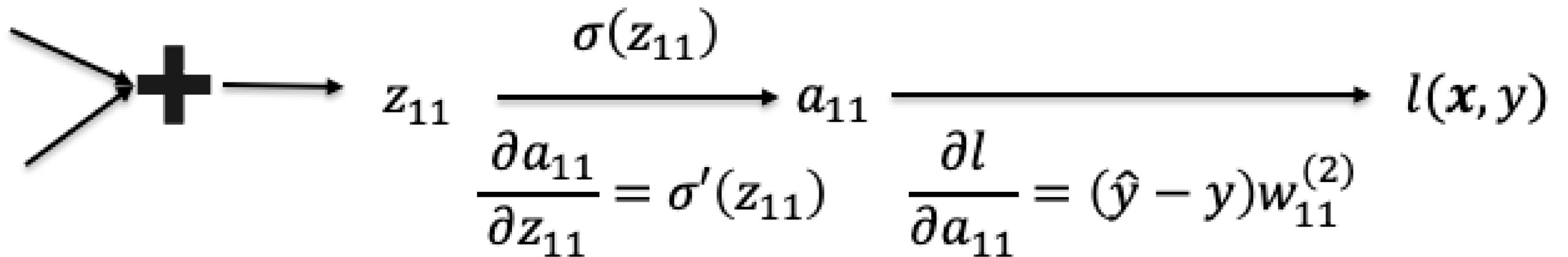
- By chain rule: $\frac{\partial l}{\partial a_{12}} = \frac{\partial l}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial a_{12}} = (\hat{y} - y) w_{12}^{(2)}$

Calculate Gradient (on one data point)



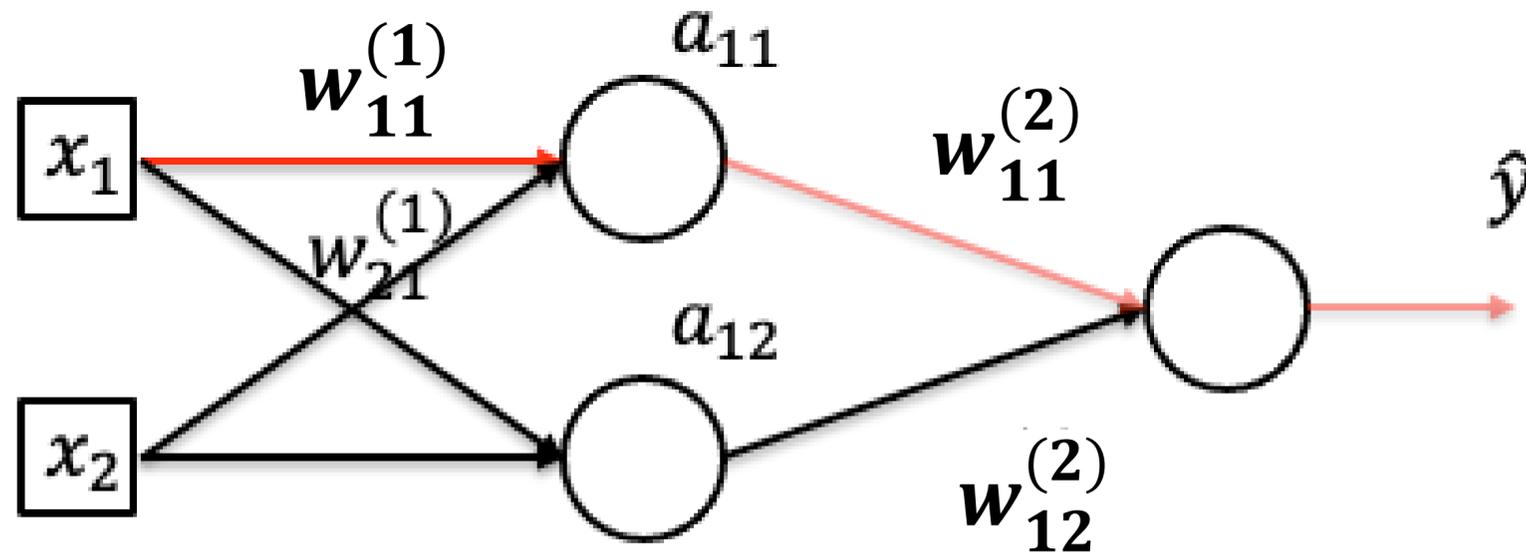
$$w_{11}^{(1)} x_1$$

$$w_{12}^{(2)} x_2$$



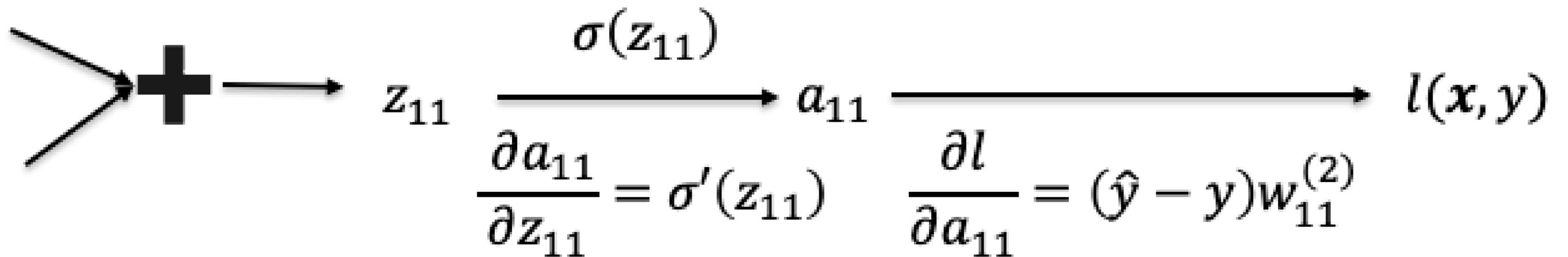
- By chain rule:
$$\frac{\partial l}{\partial w_{11}^{(1)}} = \frac{\partial l}{\partial a_{11}} \frac{\partial a_{11}}{\partial z_{11}} \frac{\partial z_{11}}{\partial w_{11}^{(1)}} = (\hat{y} - y)w_{11}^{(2)} a_{11}(1 - a_{11}) \frac{\partial z_{11}}{\partial w_{11}^{(1)}}$$

Calculate Gradient (on one data point)



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

$$w_{12}^{(2)} x_2$$



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

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

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 limited but enough accuracy is needed?

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

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 limited but enough accuracy is needed?

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

What we've learned today...

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

Suggested reading

- Textbook: Artificial Intelligence: A Modern Approach (4th edition). Stuart Russell and Peter Norvig. Pearson, 2020.
 - Sections 21 Introduction, 21.1, 21.2, 21.4 ,21.5