



CS 540 Introduction to Artificial Intelligence

Neural Networks (II)

University of Wisconsin-Madison

Spring 2023

Announcements

- **Homeworks:**
 - HW 6 Due Tuesday, March 28
- **Midterms are being graded**
- **Class roadmap:**

Tuesday, Mar 21	ML: Neural Networks II
Thursday, Mar 23	ML: Neural Networks III
Tuesday, Mar 28	Deep Learning I
Thursday, Mar 30	Deep Learning II

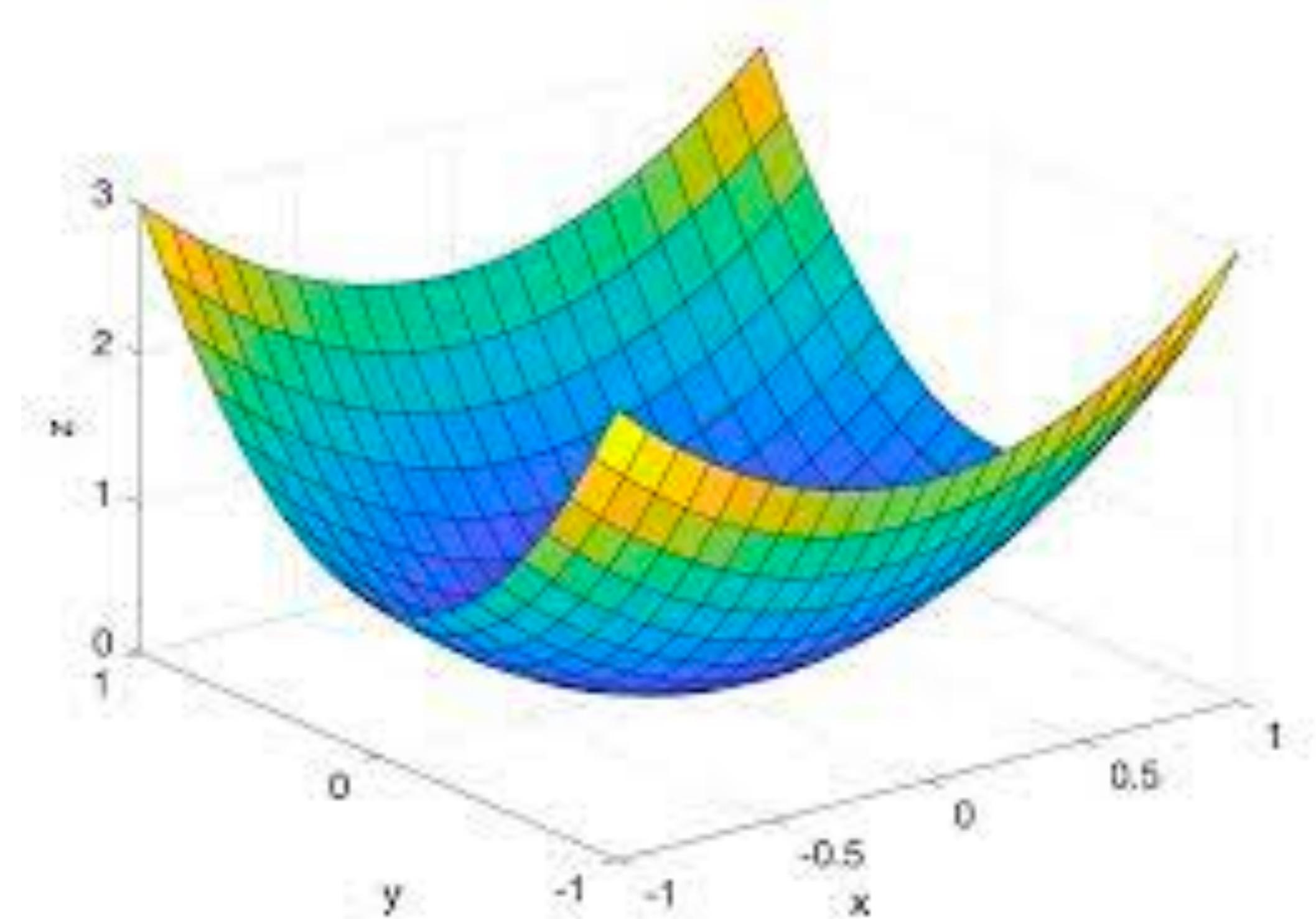
Section 001 Mid-Semester Evaluation

- Lecture
 - Slides not matching website?
 - Amount of material in some lectures.
- Assignments
 - Feedback on what is going wrong.
- Exams
 - Wanting to see more practice problems before the exam.

Today's outline

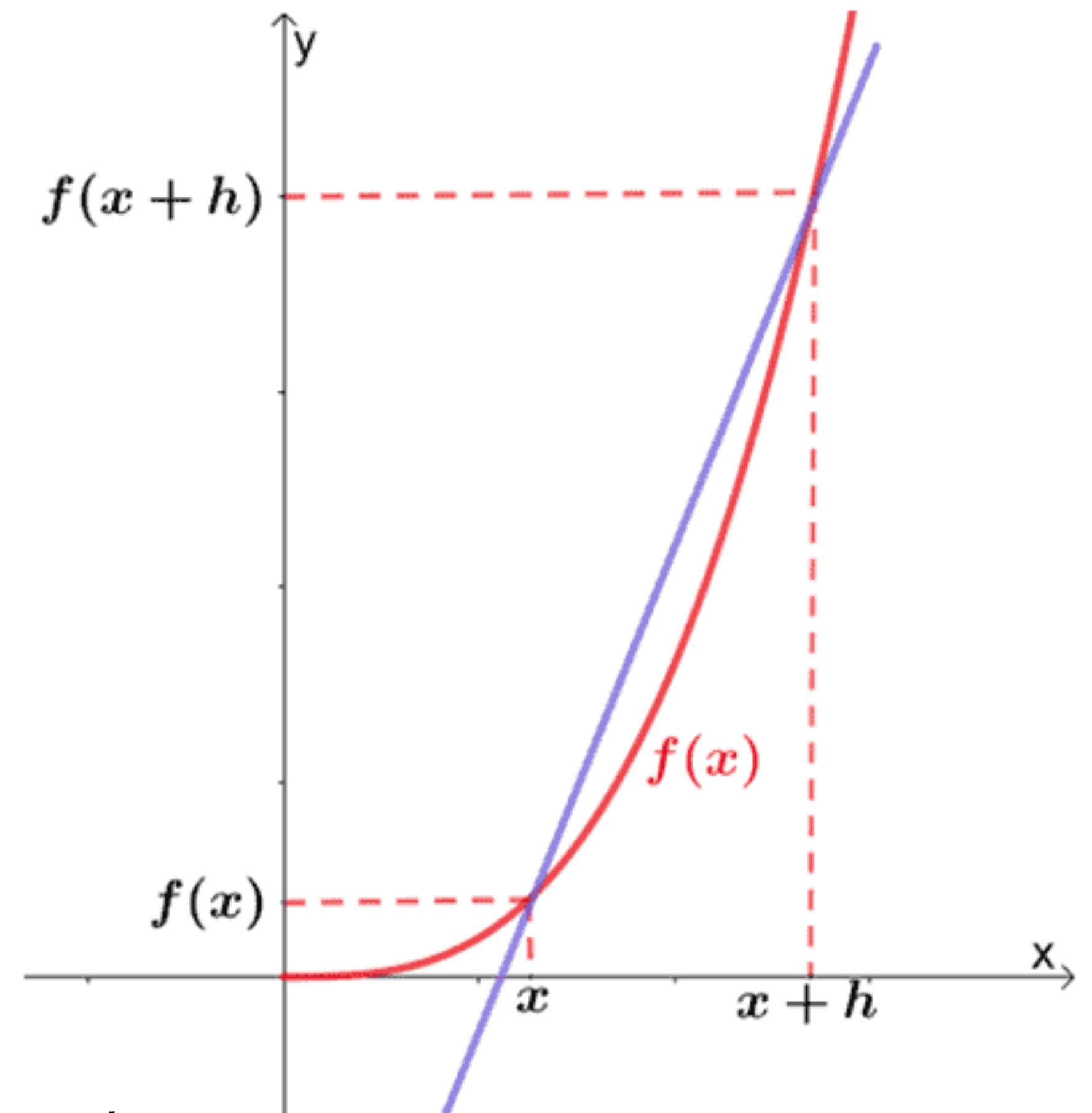
- Multivariate Calculus Intro / Review
- Review of Multi-layer Perceptron
 - Single output
 - Multiple output
- How to train neural networks
 - Gradient descent
 - Backpropagation

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 $\frac{df}{dx} = cpx^{p-1}$

- 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 $f(x) = f_1(f_2(\dots f_n(x)\dots))$ then $\frac{df}{dx} = \frac{df_1}{df_2} \frac{df_2}{df_3} \dots \frac{df_n}{dx}$

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.
 - $\frac{df}{d\mathbf{x}} = [\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n}]$

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

- 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: $\frac{\partial \sigma}{\partial z} = \sigma(z)(1 - \sigma(z))$.

Quiz Break

- 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: $\frac{\partial \sigma}{\partial z} = \sigma(z)(1 - \sigma(z))$.

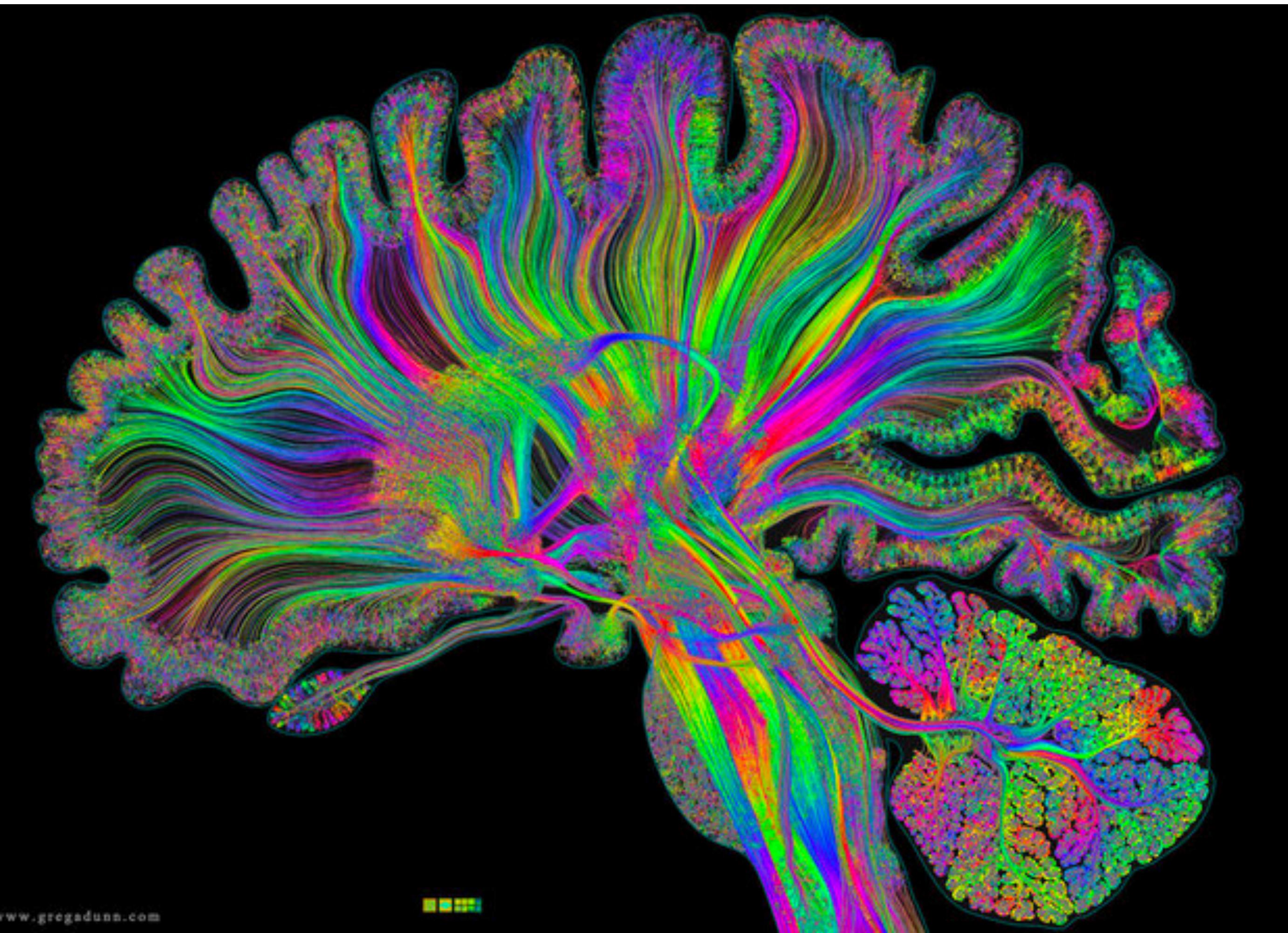
Let $a = \sigma(b)$

Let $b = w_1 x_1 + w_2 x_2$

$$\frac{\partial f}{\partial w_1} = \frac{\partial f}{\partial a} \frac{\partial a}{\partial b} \frac{\partial b}{\partial w_1}$$

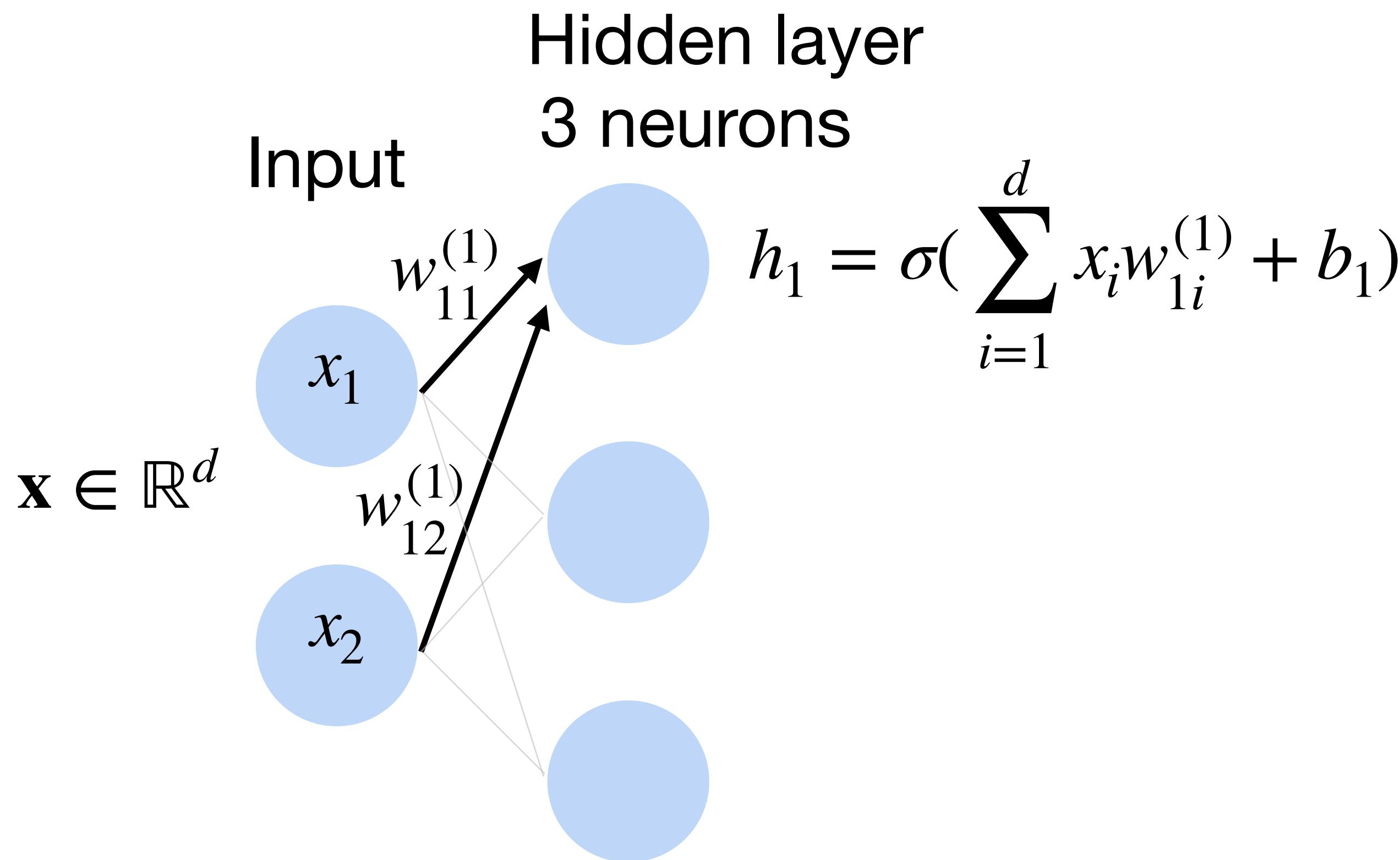
$$\frac{\partial f}{\partial w_1} = \frac{y}{a} \sigma(b)(1 - \sigma(b))x_2 = (1 - \sigma(w_1 x_1 + w_2 x_2))x_2$$

Multilayer Perceptron



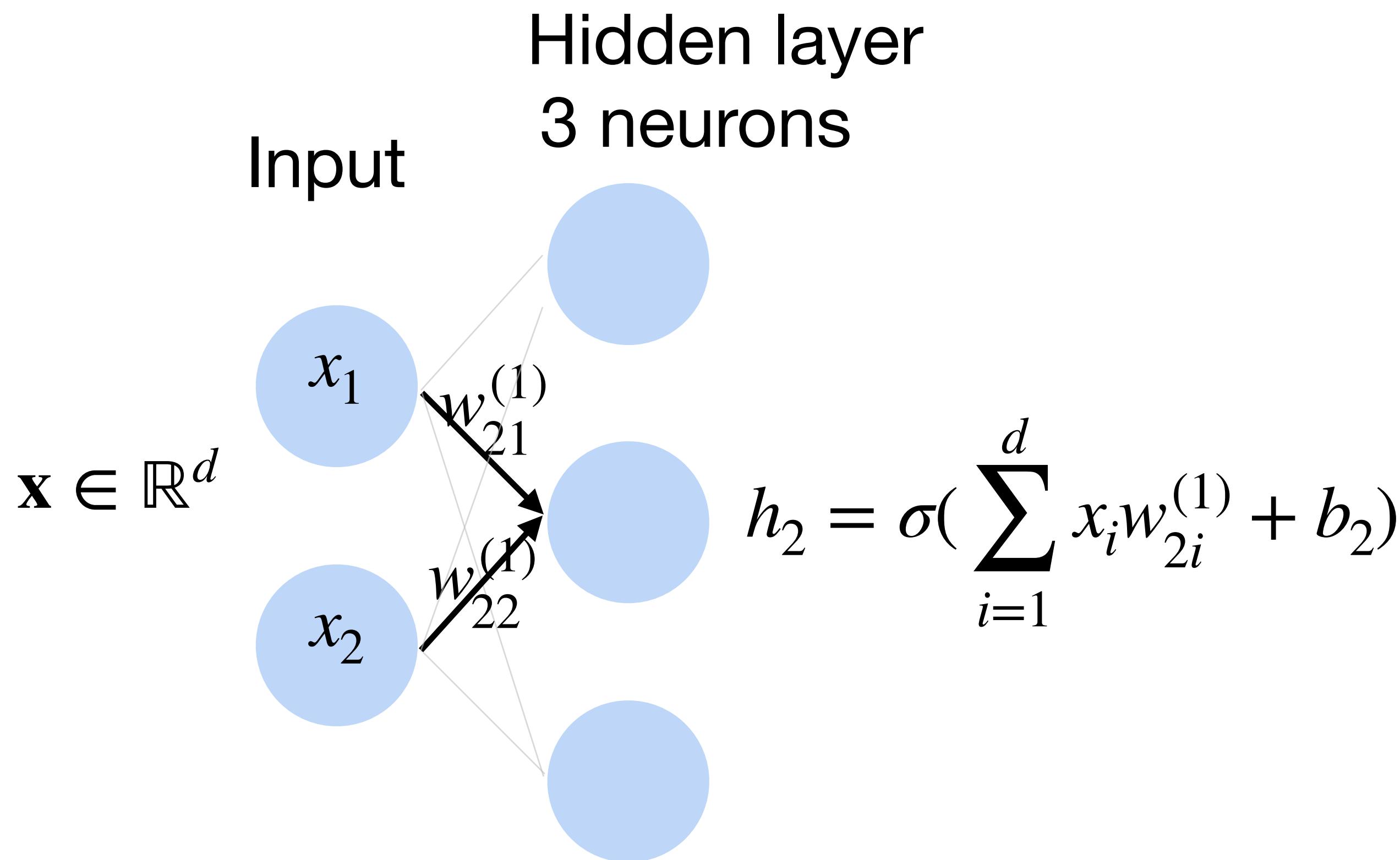
Multi-layer perceptron: Example

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



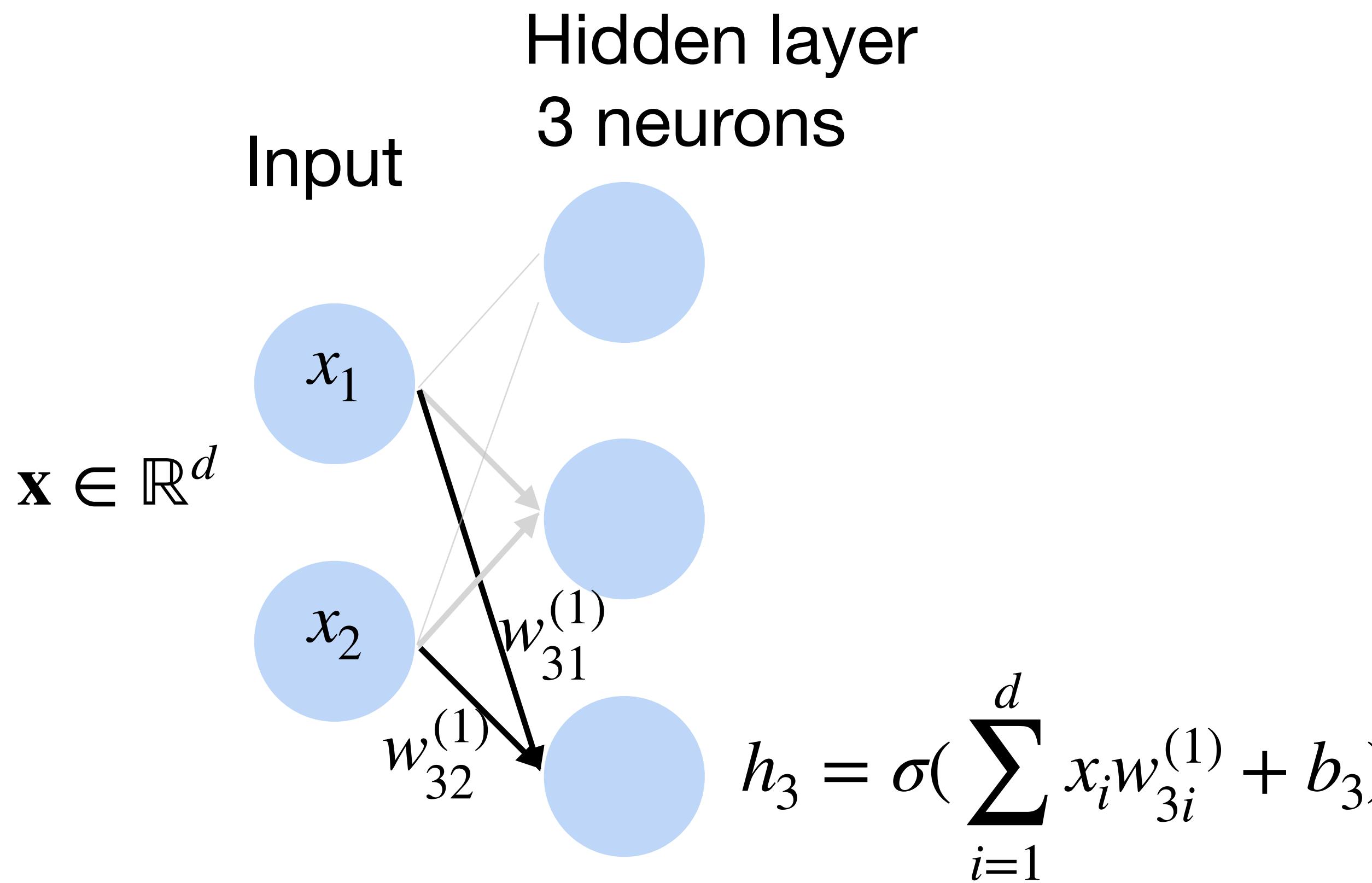
Multi-layer perceptron: Example

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



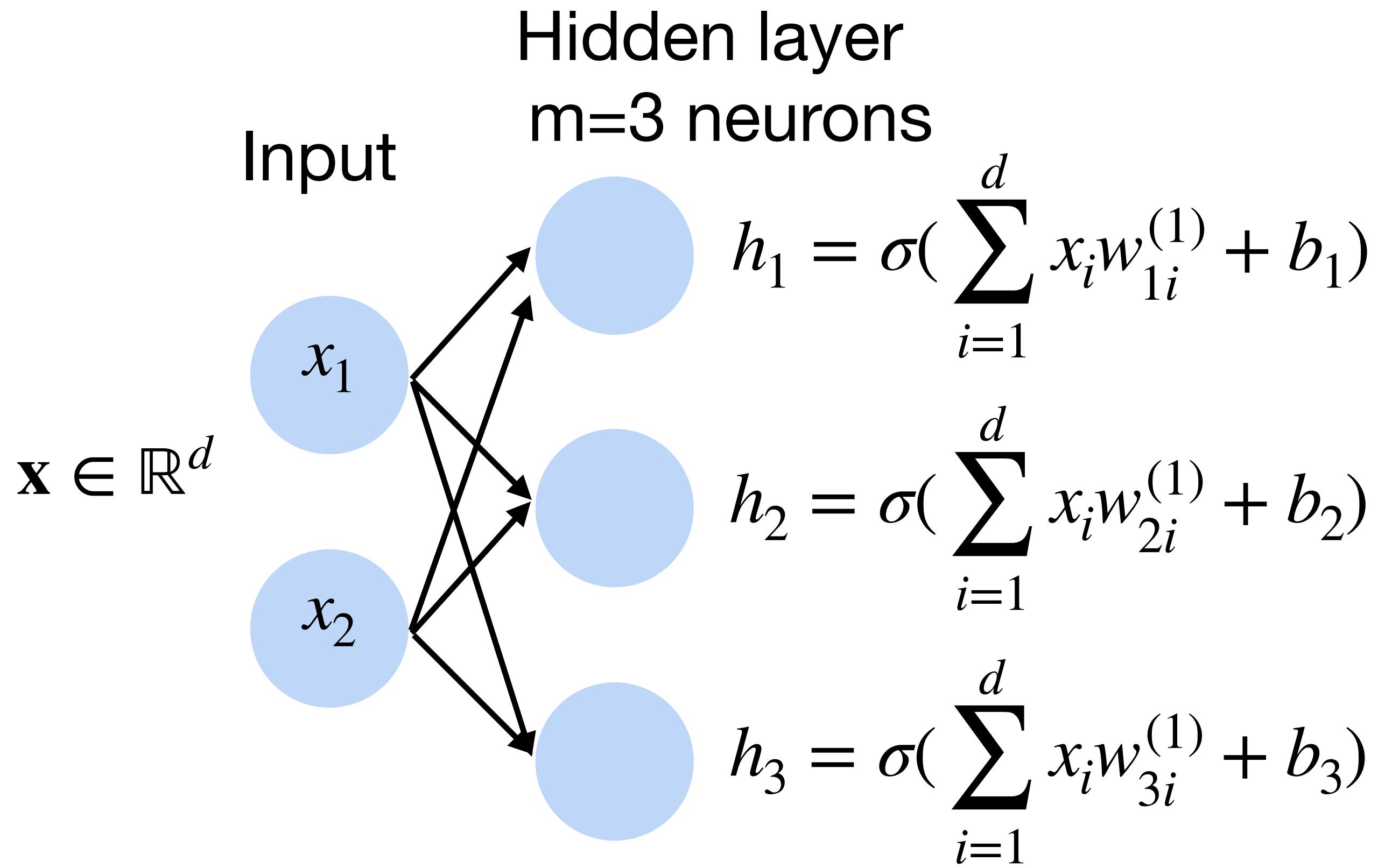
Multi-layer perceptron: Example

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



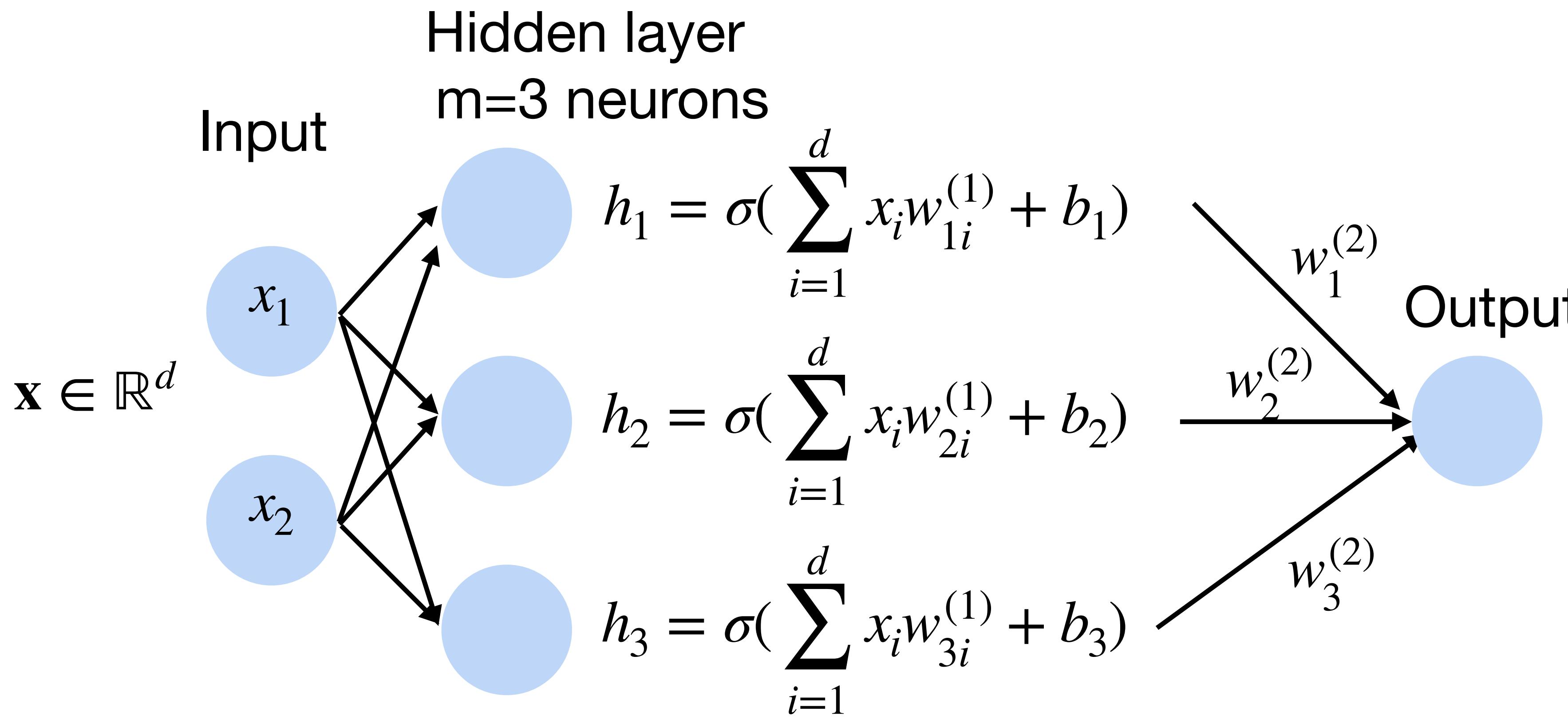
Multi-layer perceptron: Example

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



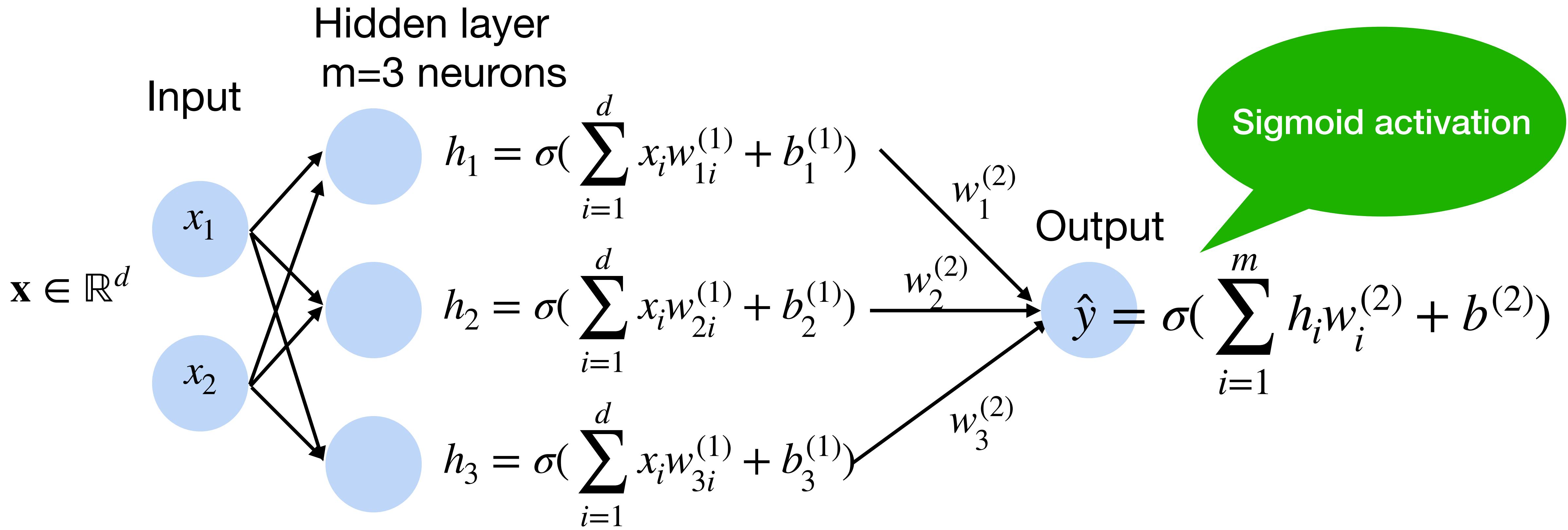
Multi-layer perceptron: Example

- 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

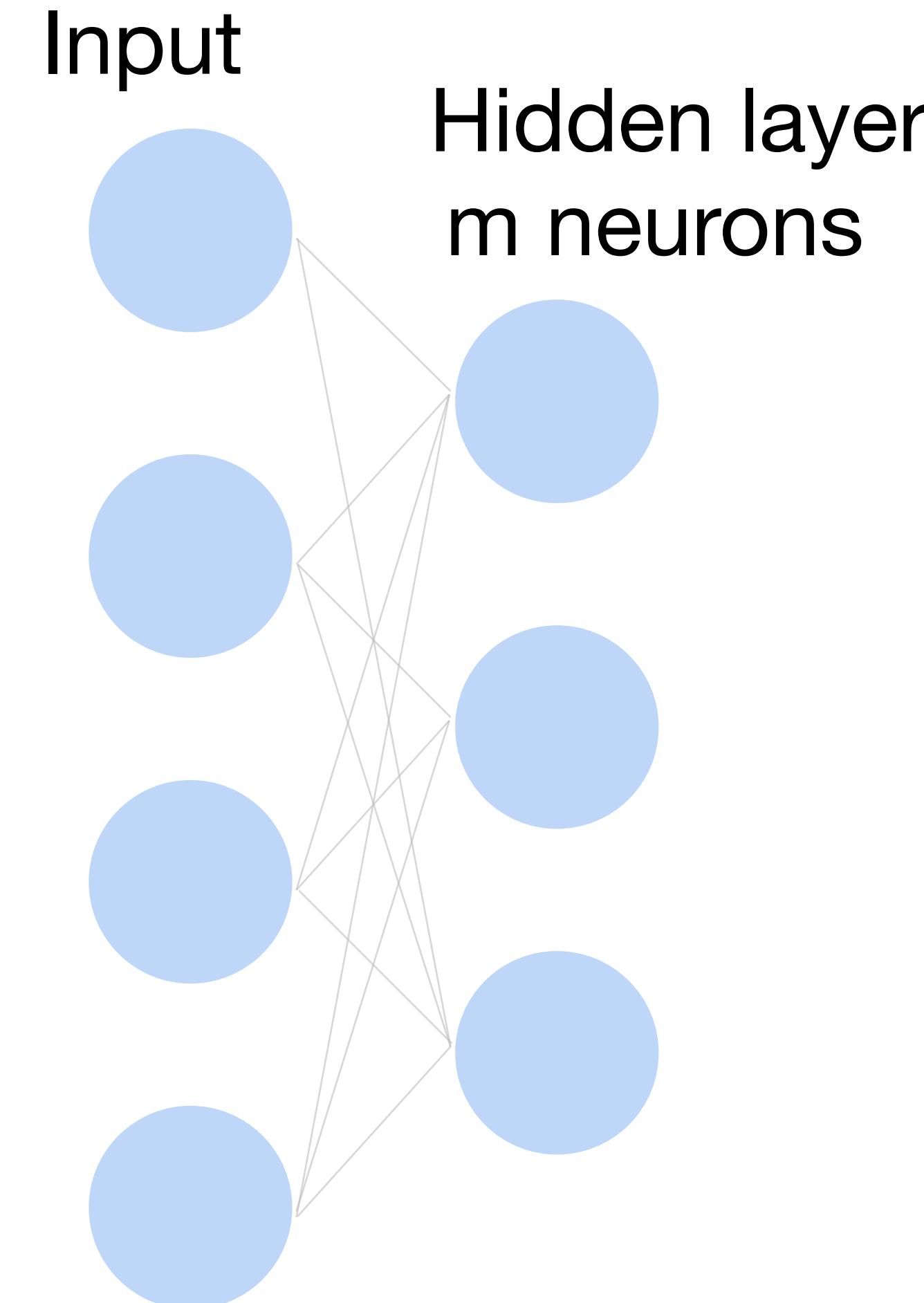


Multi-layer perceptron: Matrix Notation

- 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

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

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

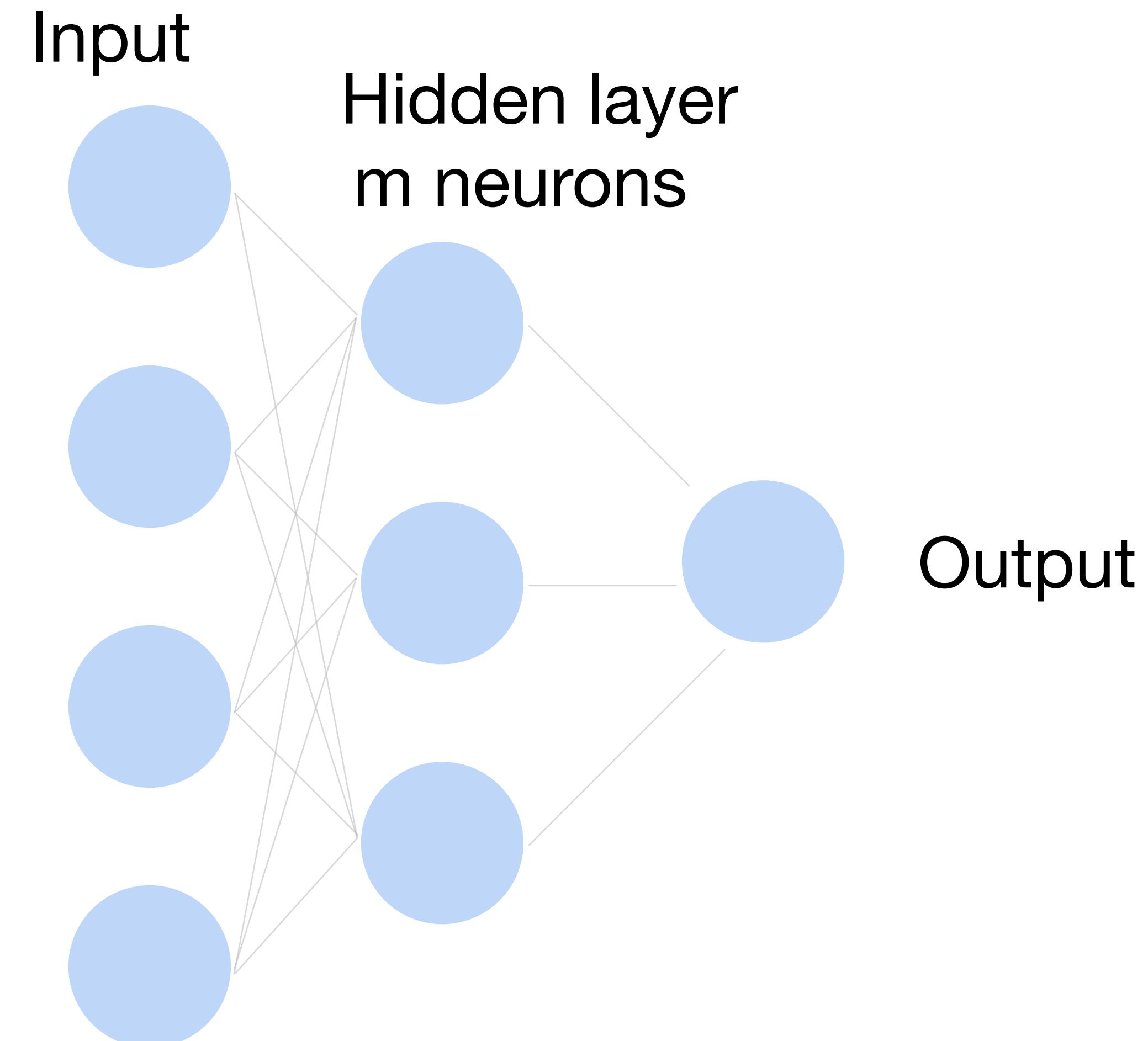


Multi-layer perceptron: Matrix Notation

- 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

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

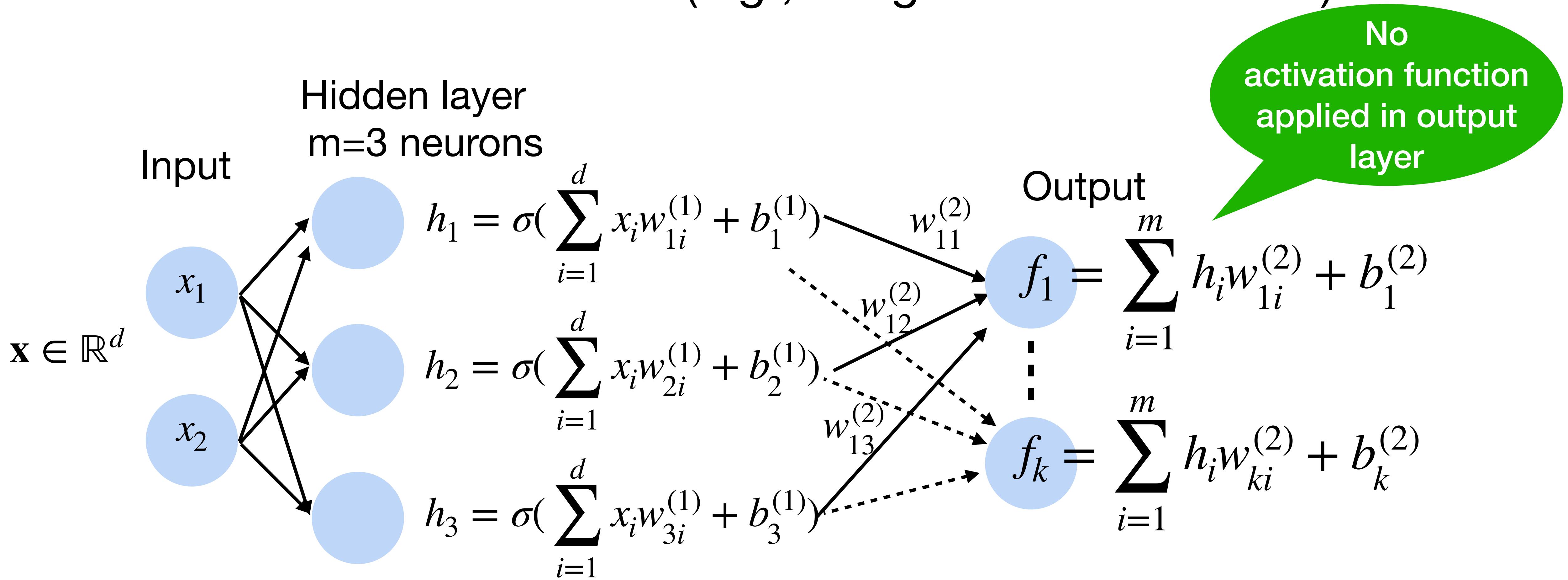
$$\hat{y} = \sigma(\mathbf{w}^{(2)}\mathbf{h} + \mathbf{b}^{(2)})$$



Neural network for K-way classification

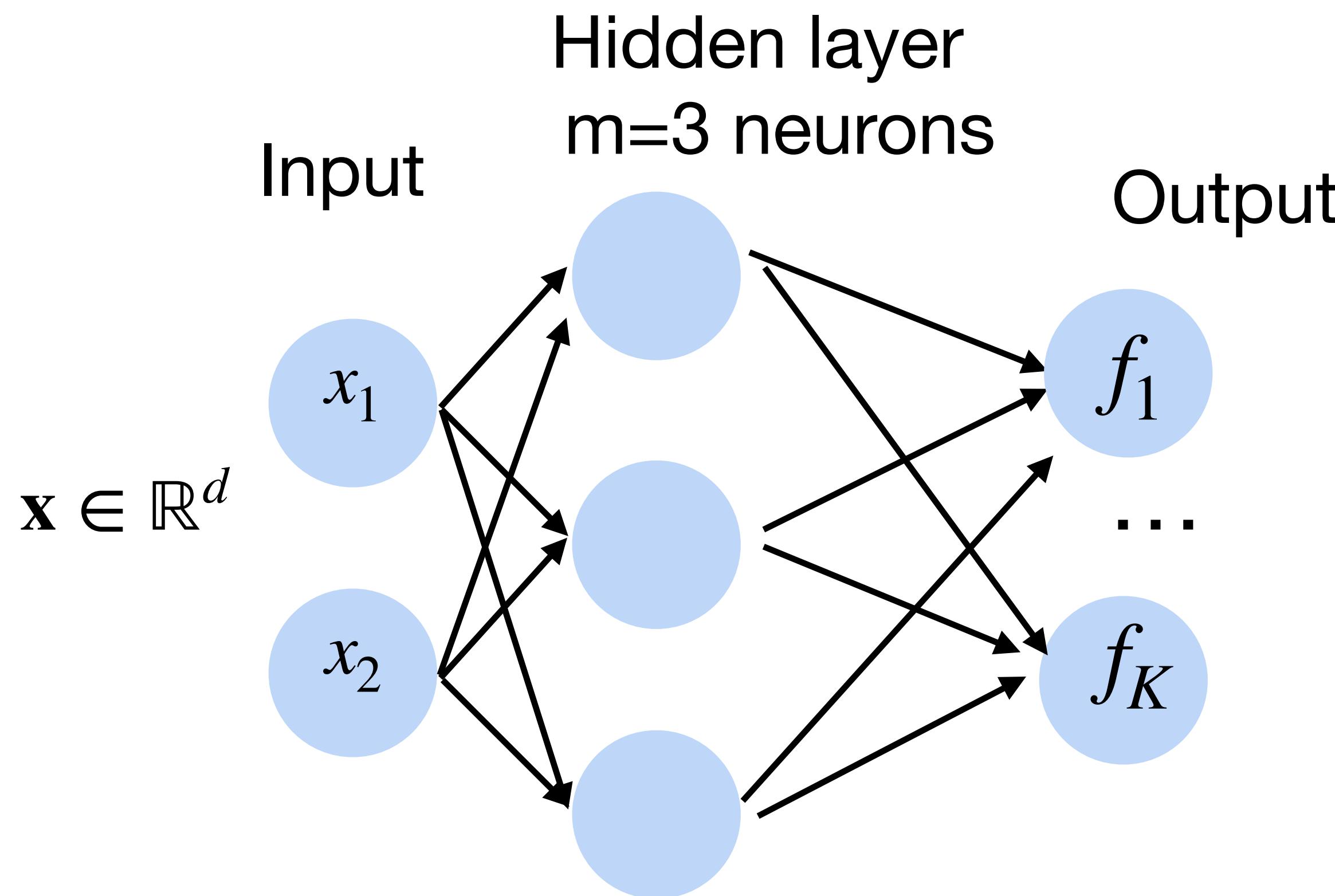
- K outputs in the final layer

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



Softmax

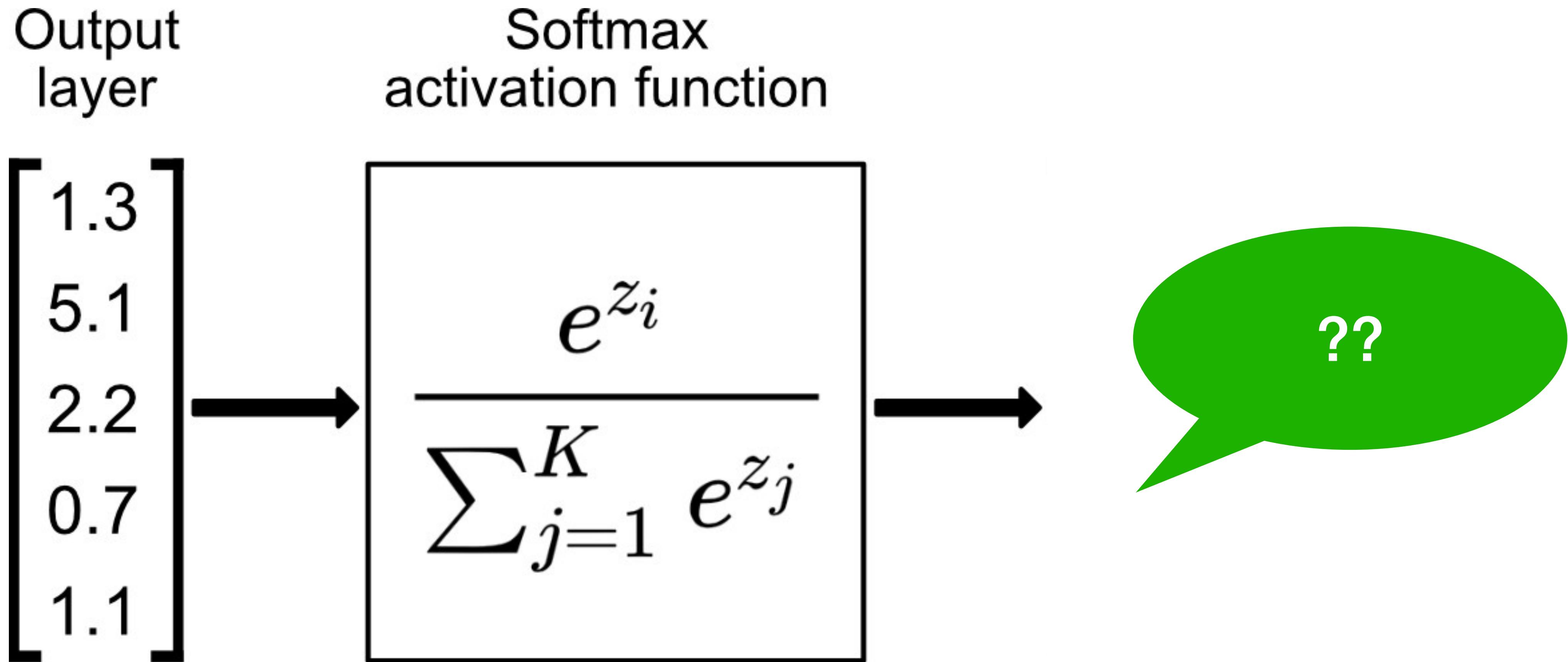
Turns outputs f into probabilities (sum up to 1 across K classes)



$$p(y | \mathbf{x}) = \text{softmax}(f) = \frac{\exp(f_y(x))}{\sum_{k=1}^K \exp(f_k(x))}$$

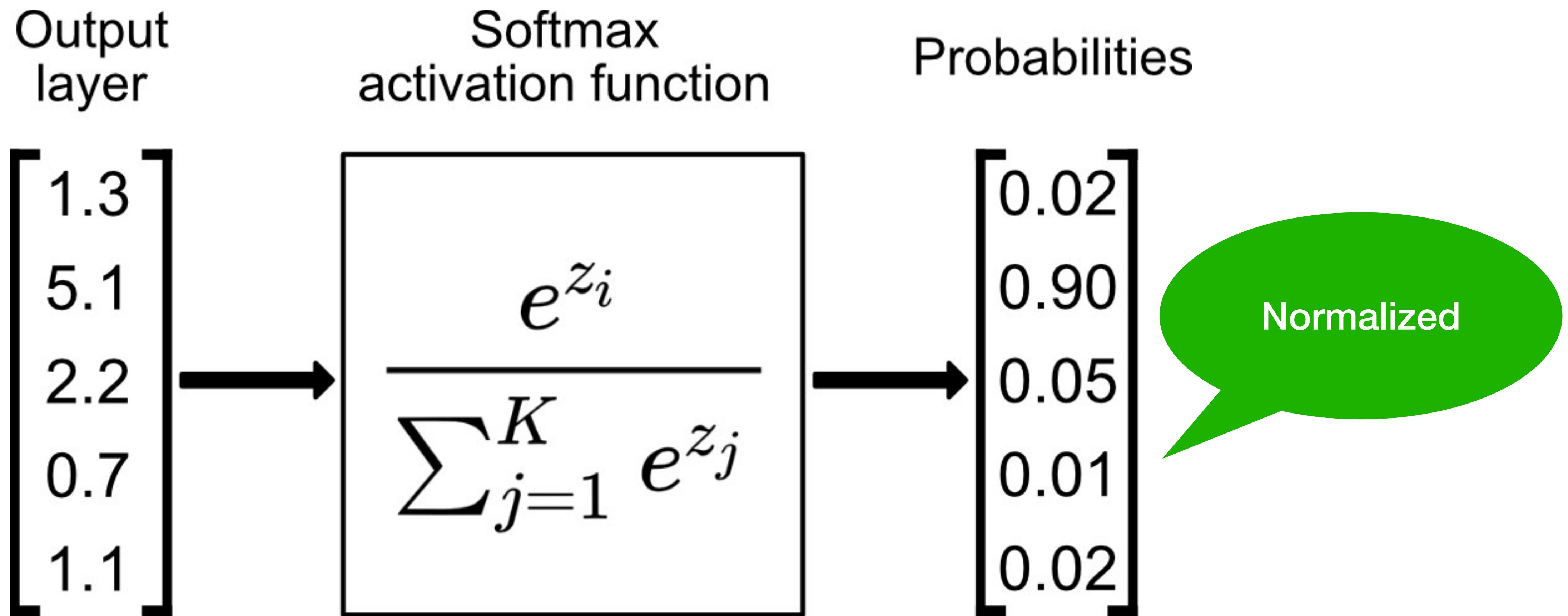
Softmax

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



Softmax

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



More complicated neural networks: multiple hidden layers

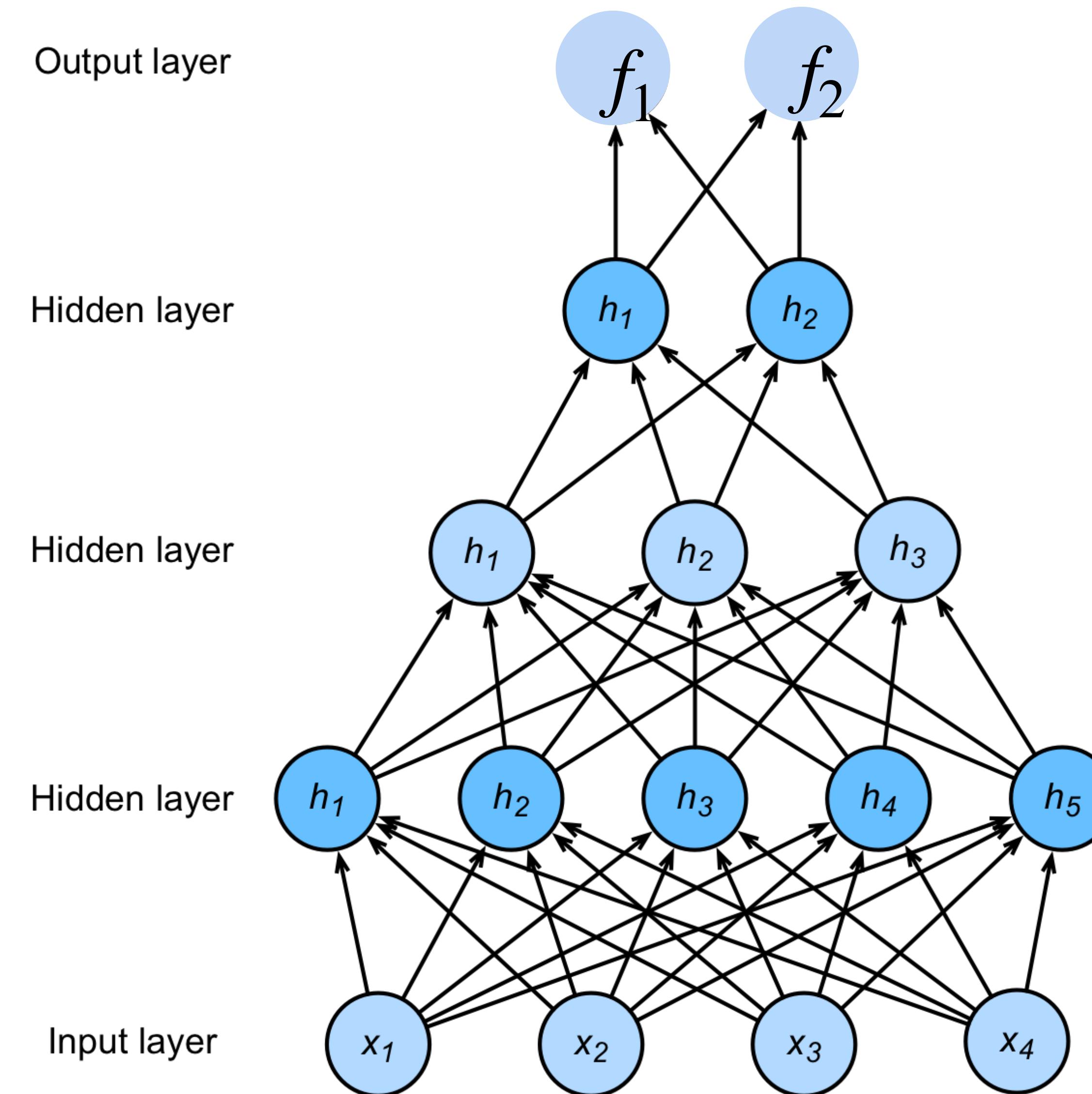
$$\mathbf{h}_1 = \sigma(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)})$$

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

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

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

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



Quiz Break

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

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

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?

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

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

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

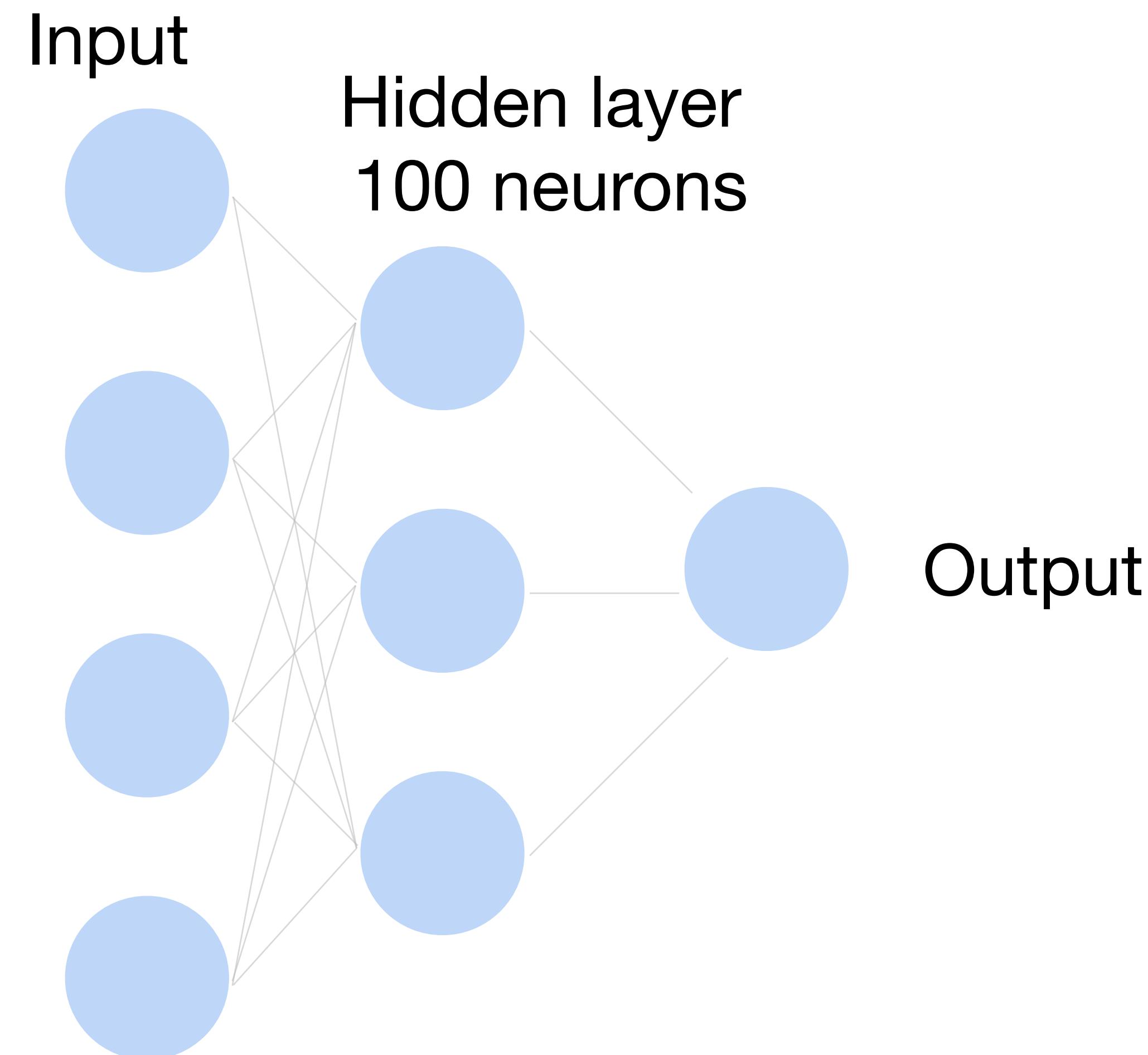


Solution:

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

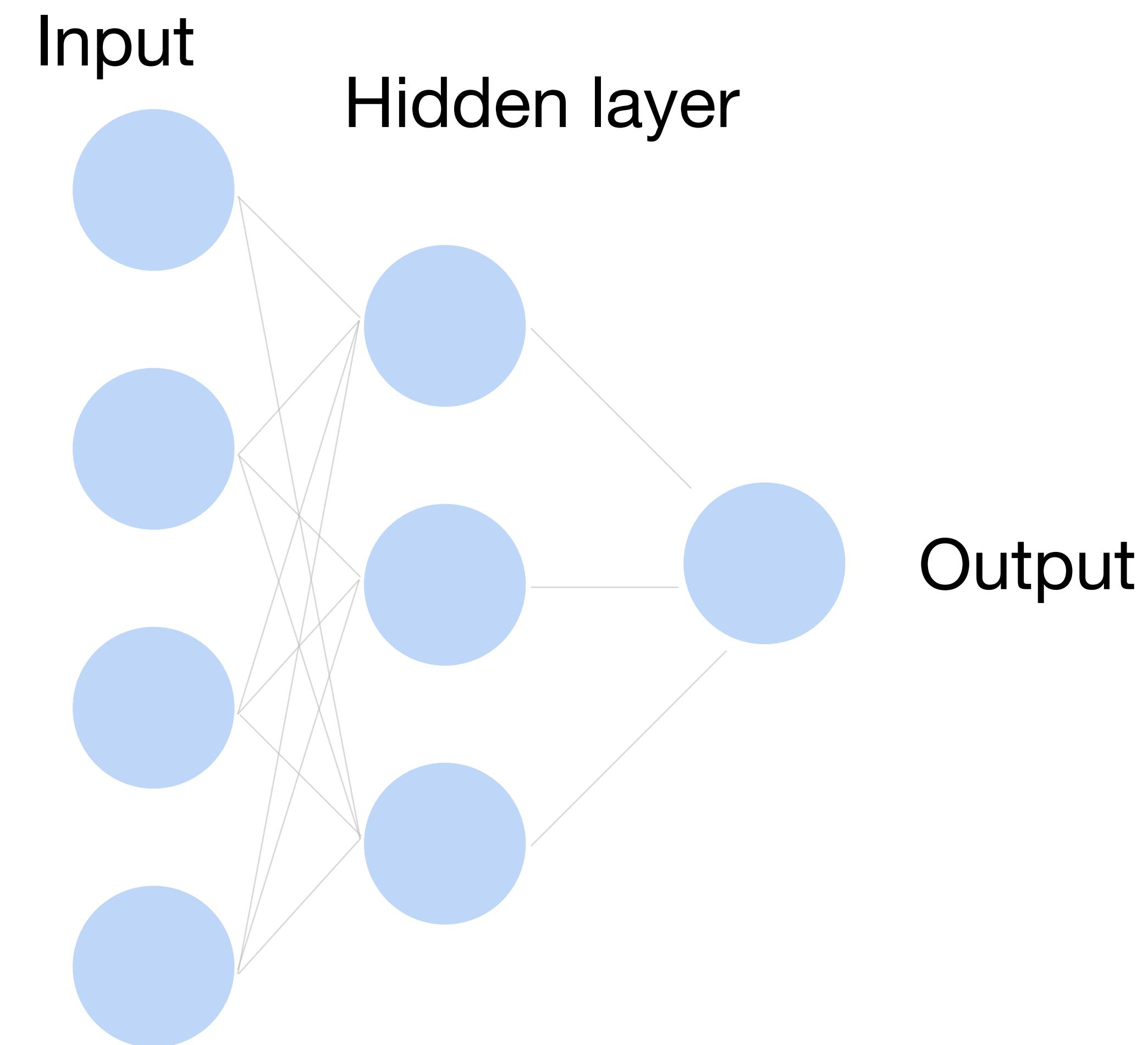
$\hat{y} \in [0,1]$ Model output for example \mathbf{x}

(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 $y = 0$

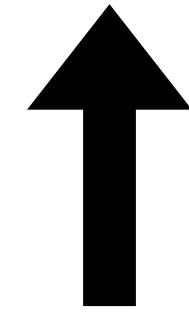


How to train a neural network? Binary classification

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

Per-sample loss:

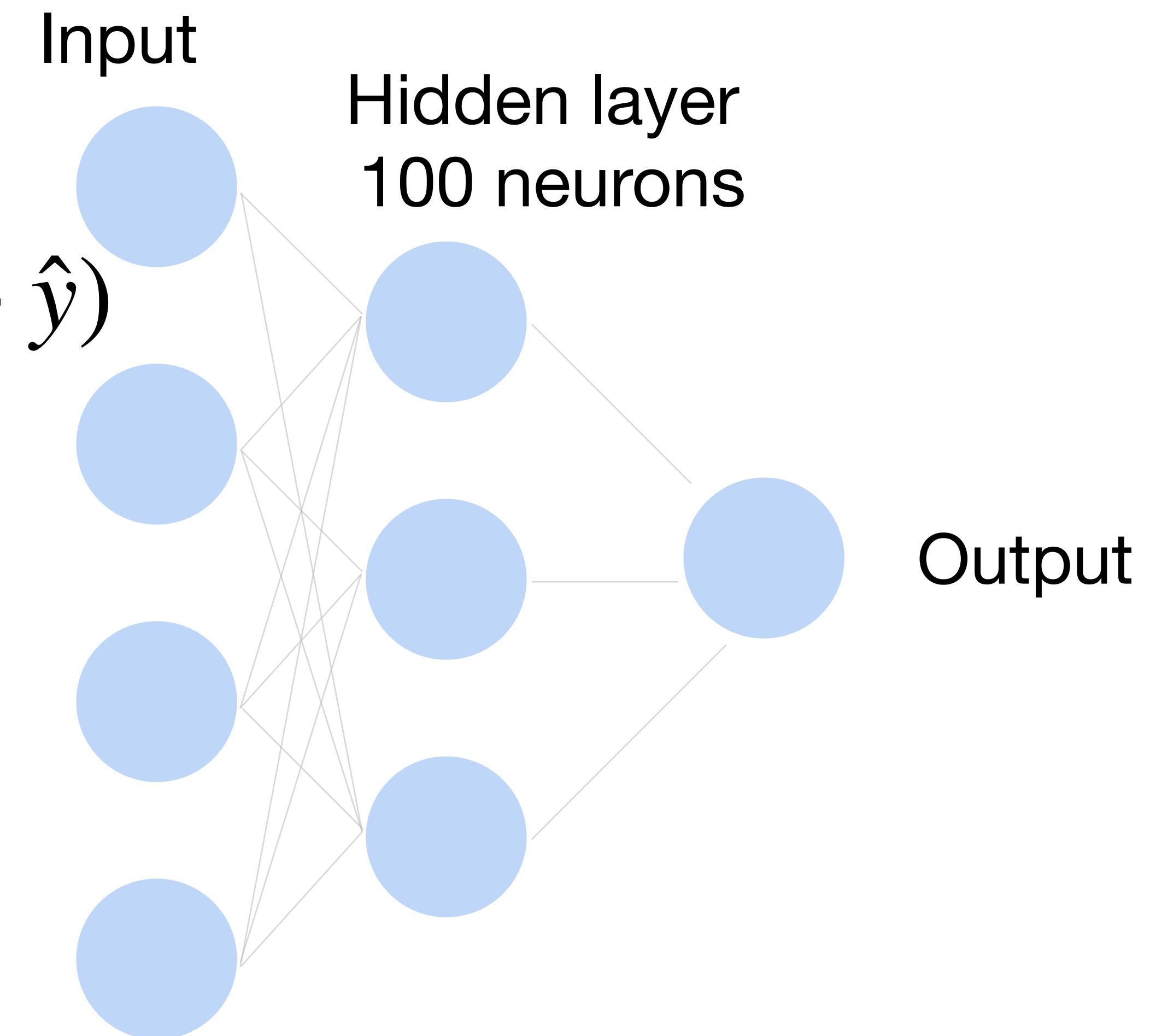
$$\ell(\mathbf{x}, y) = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y})$$



Negative log likelihood

Minimizing NLL is equivalent to Max Likelihood Learning (MLE)

Also known as **binary cross-entropy**



How to train a neural network? Multiclass

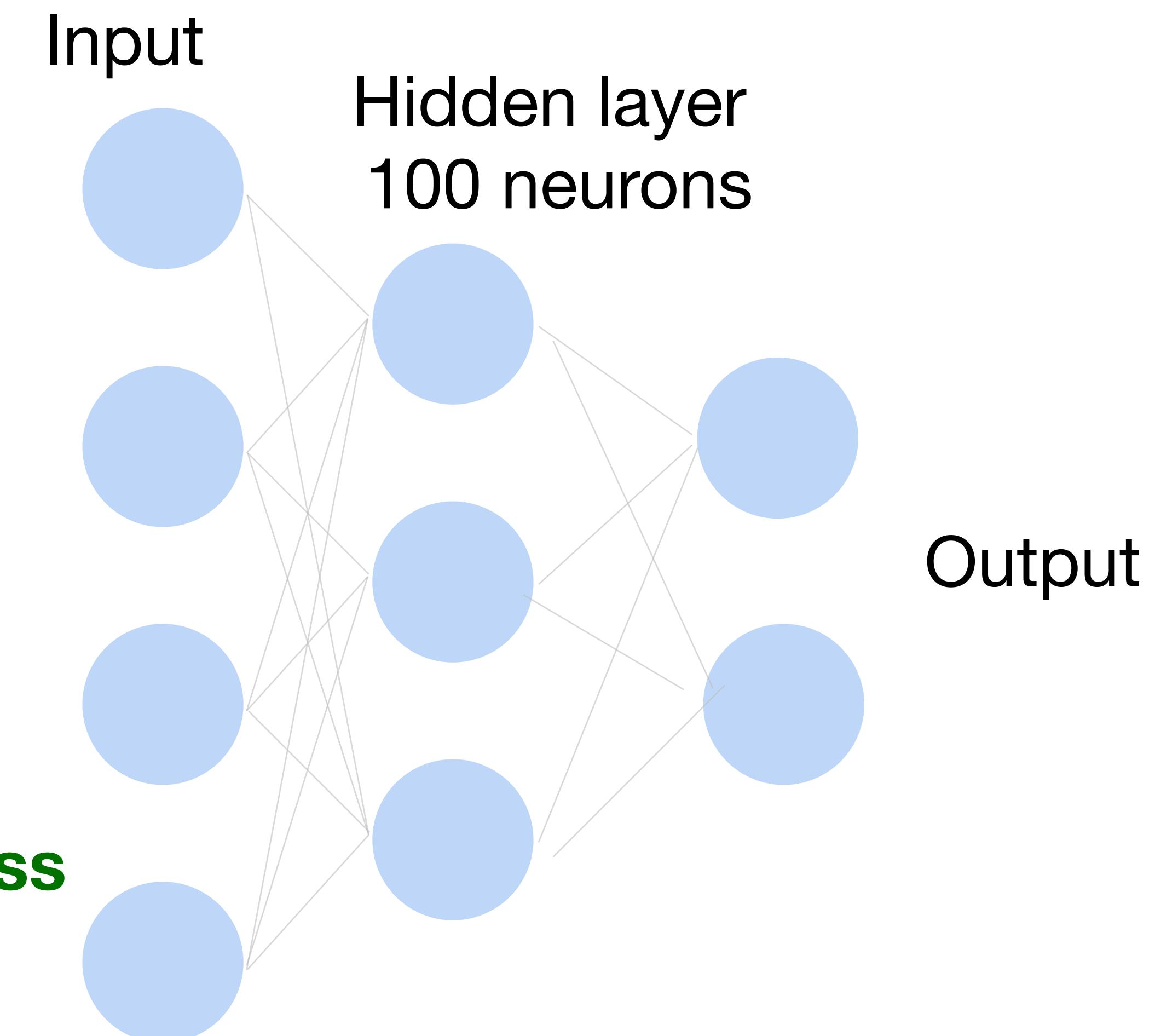
Loss function: $\frac{1}{|D|} \sum_{(\mathbf{x}, y) \in D} \ell(\mathbf{x}, y)$

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

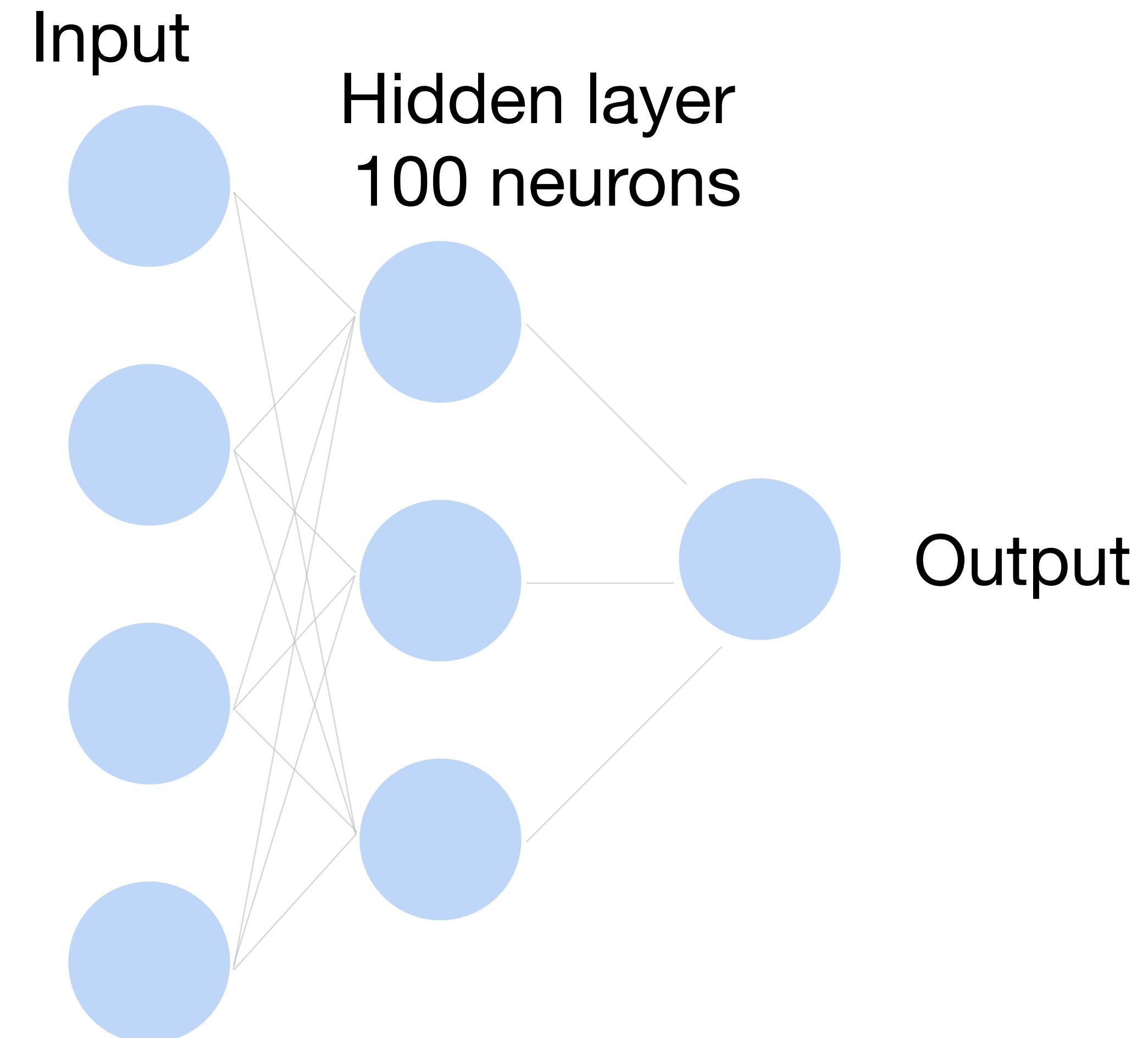


How to train a neural network?

Update the weights W to minimize the loss function

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

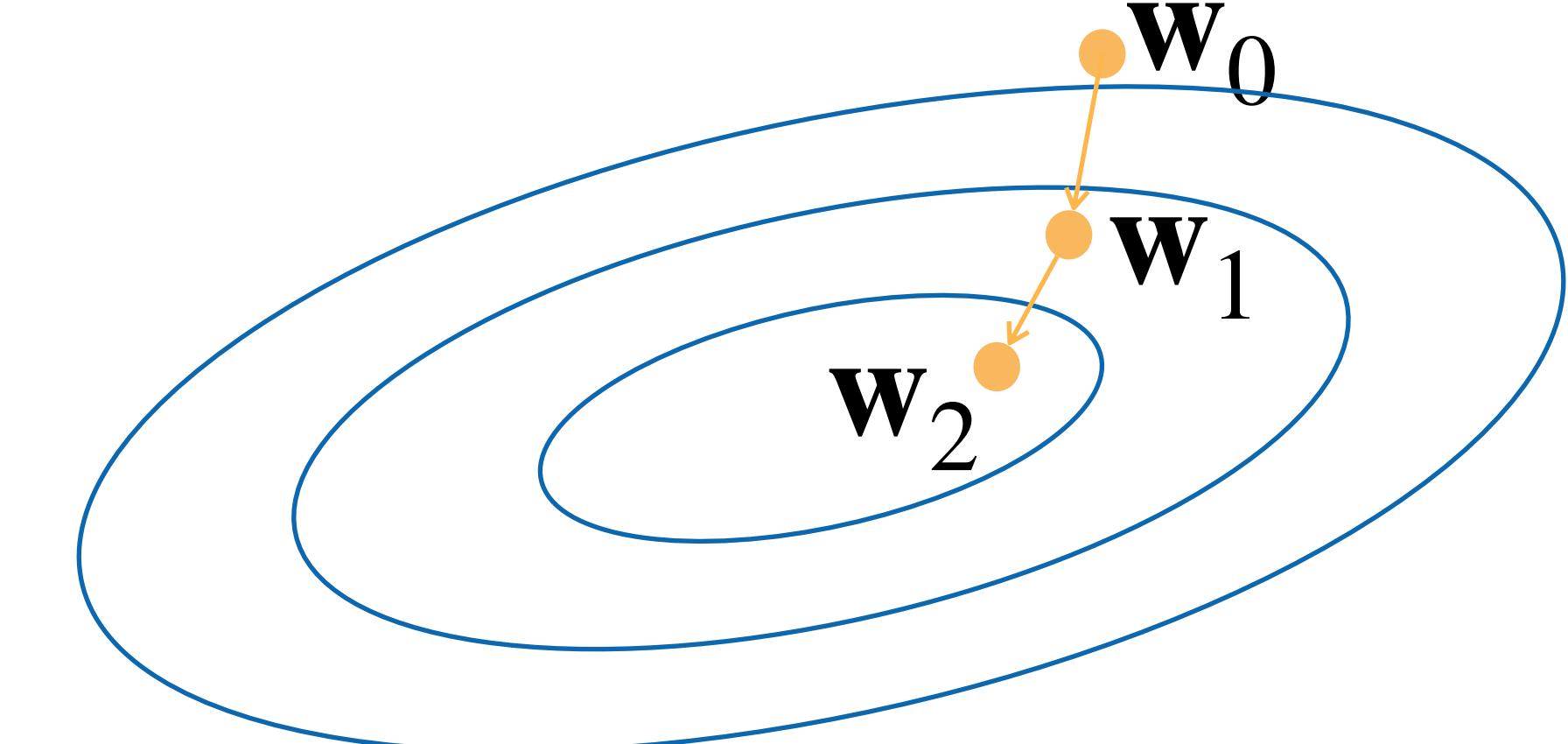
Use gradient descent!



Gradient Descent

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

$$\begin{aligned} w_t &= w_{t-1} - \alpha \frac{\partial L}{\partial w_{t-1}} \\ &= w_{t-1} - \alpha \frac{1}{|D|} \sum_{(x,y) \in D} \frac{\partial \ell(x, y)}{\partial w_{t-1}} \end{aligned}$$



D can
be very large.
Expensive

- Repeat until converges

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

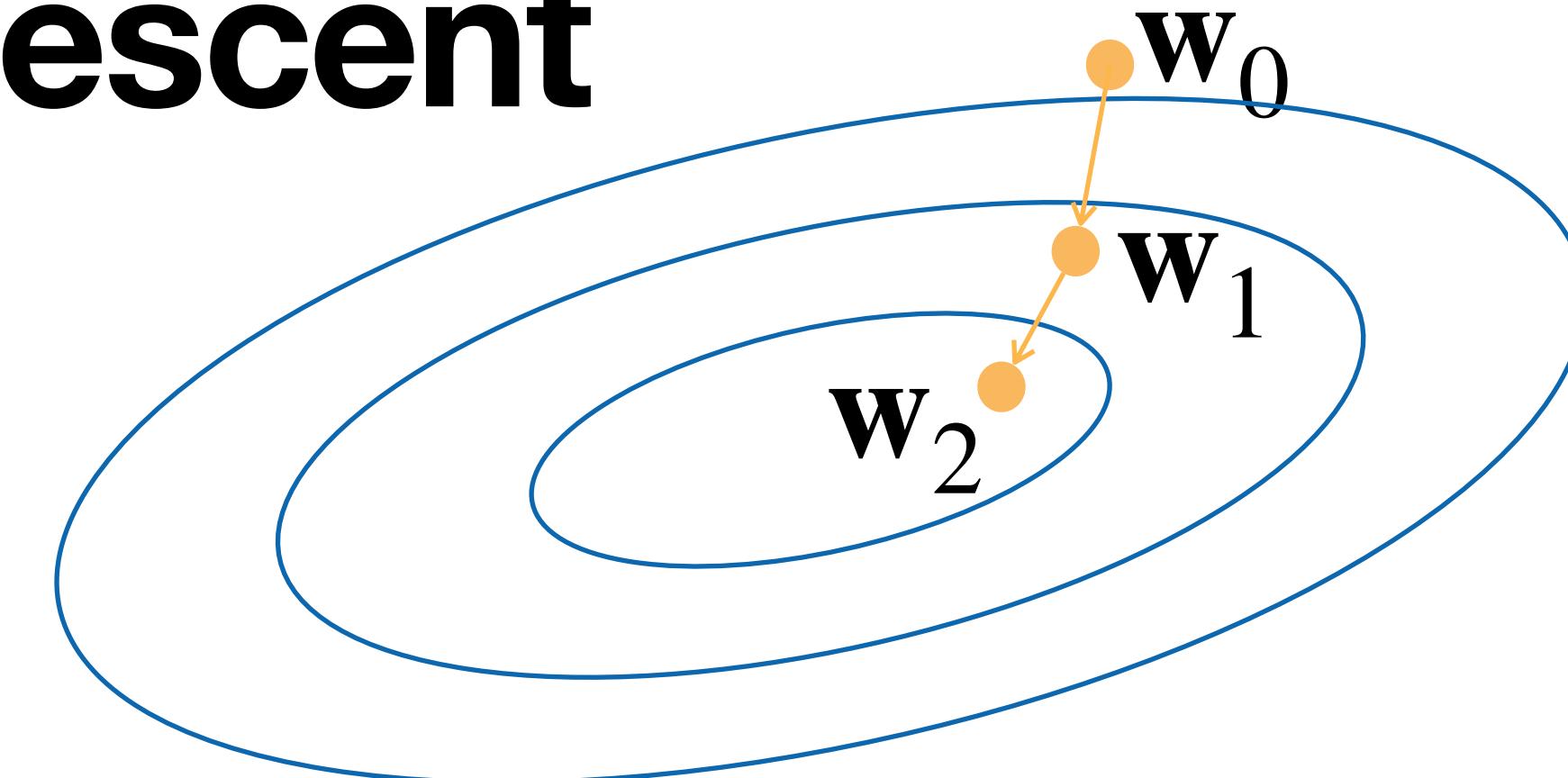
Minibatch Stochastic Gradient Descent

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

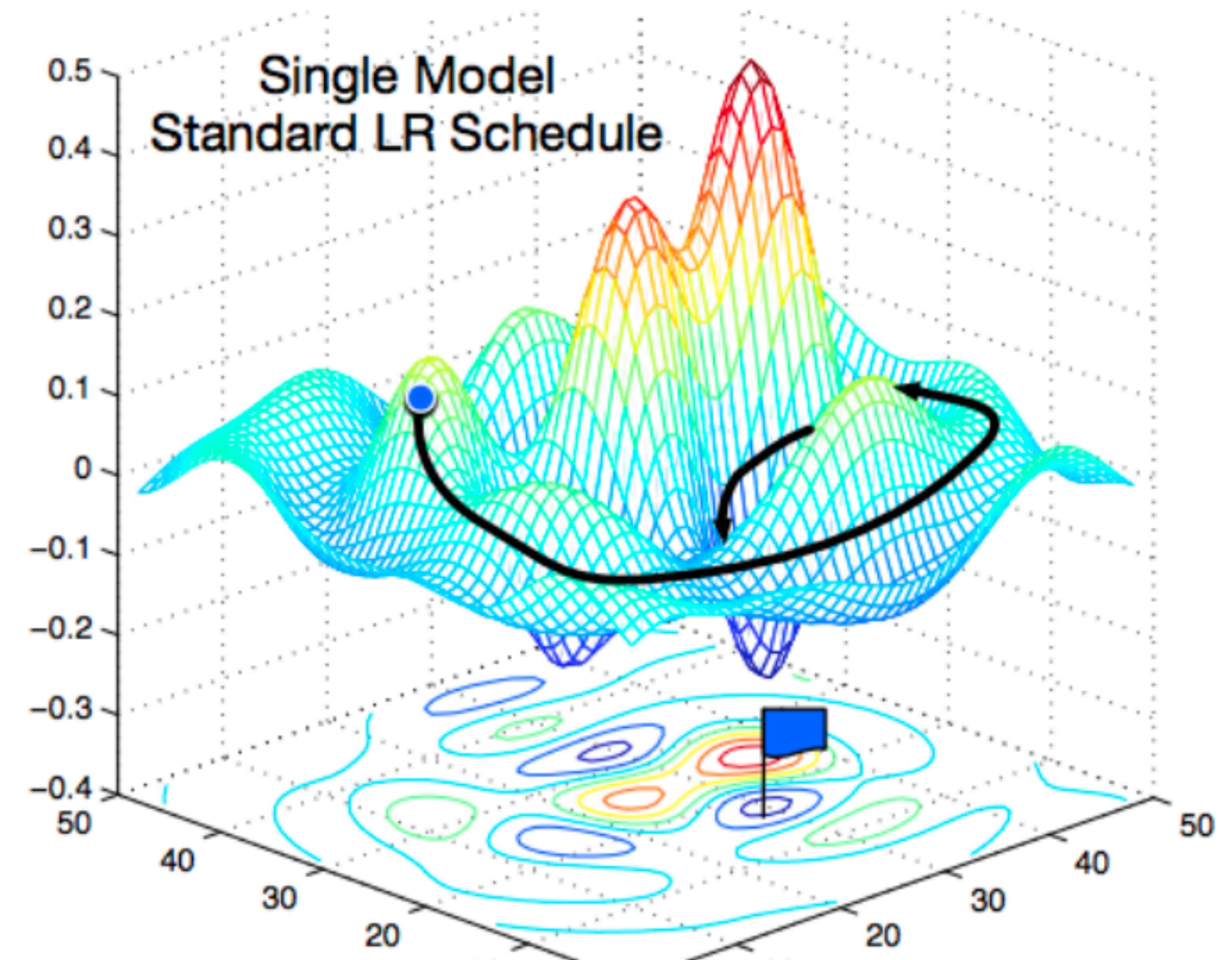
- Randomly sample a subset (mini-batch) $B \subset D$
Update parameters:

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

- Repeat until converges

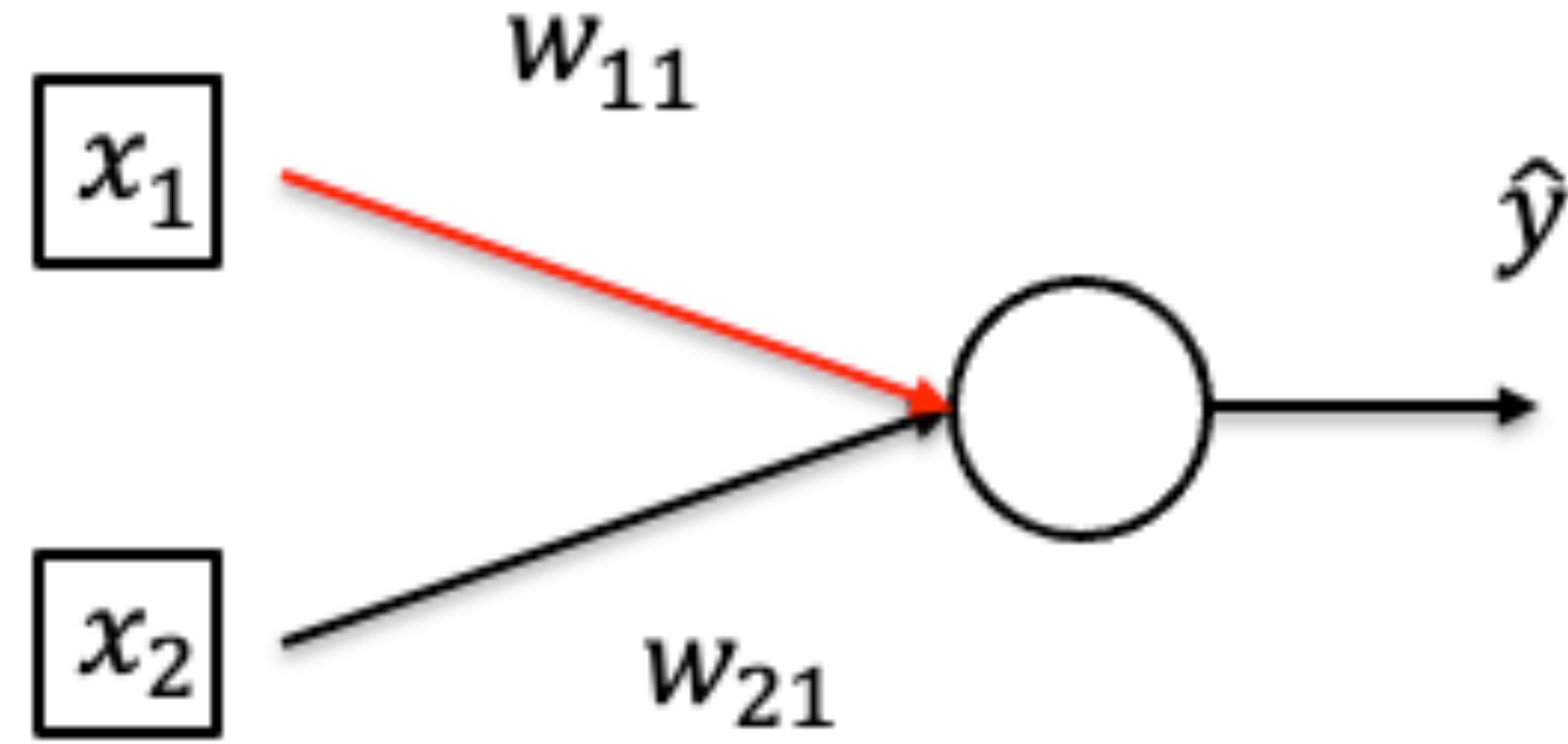


Non-convex Optimization



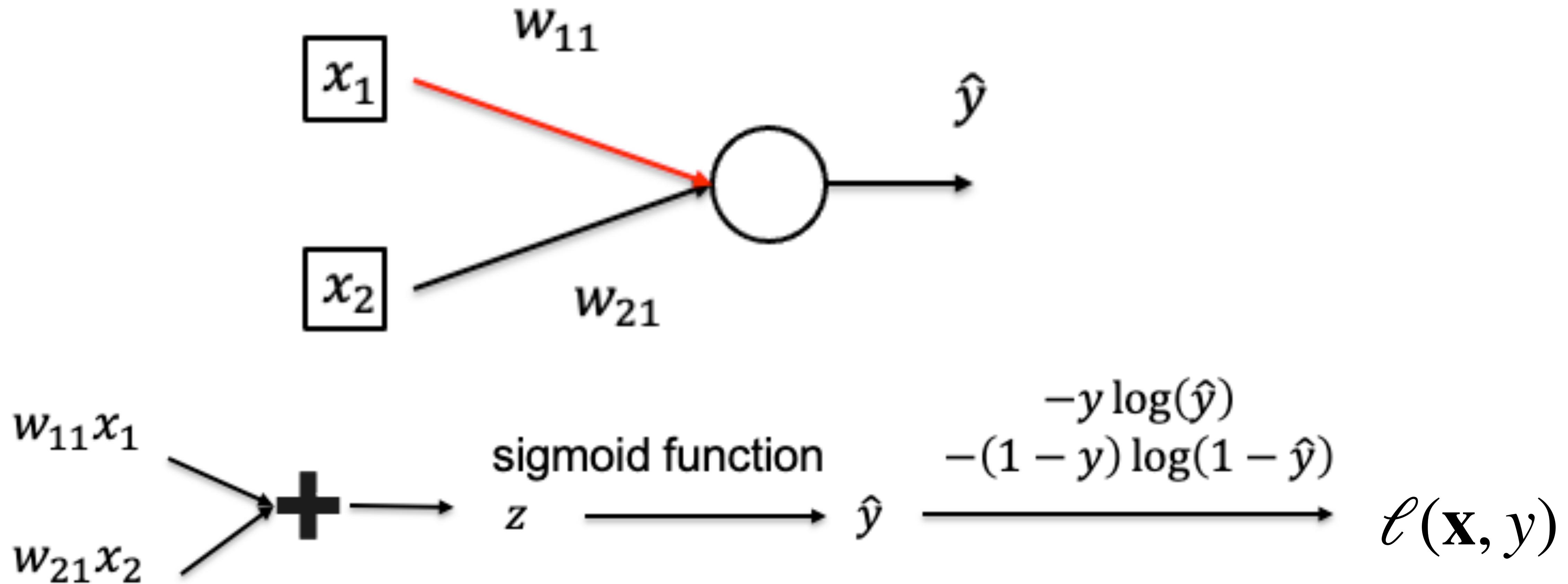
[Gao and Li et al., 2018]

Calculate Gradient (on one data point)



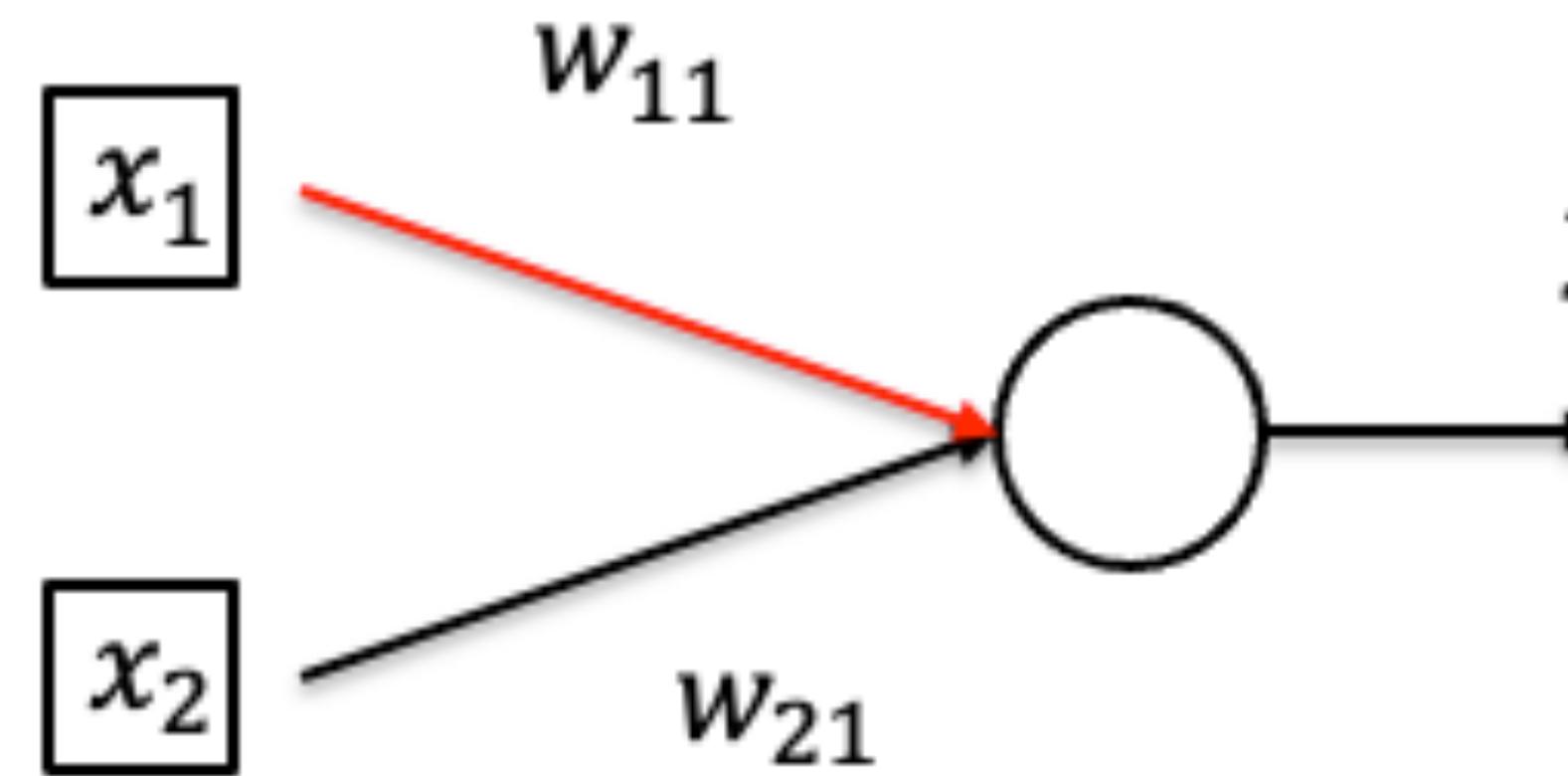
- Want to compute $\frac{\partial \ell(\mathbf{x}, y)}{\partial w_{11}}$
- Data point: $((x_1, x_2), y)$

Calculate Gradient (on one data point)



Use chain rule!

Calculate Gradient (on one data point)



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 .
- The sigmoid function is applied to z to produce the prediction \hat{y} .
- The loss function $\ell(\mathbf{x}, y)$ is calculated based on \hat{y} and the target y .

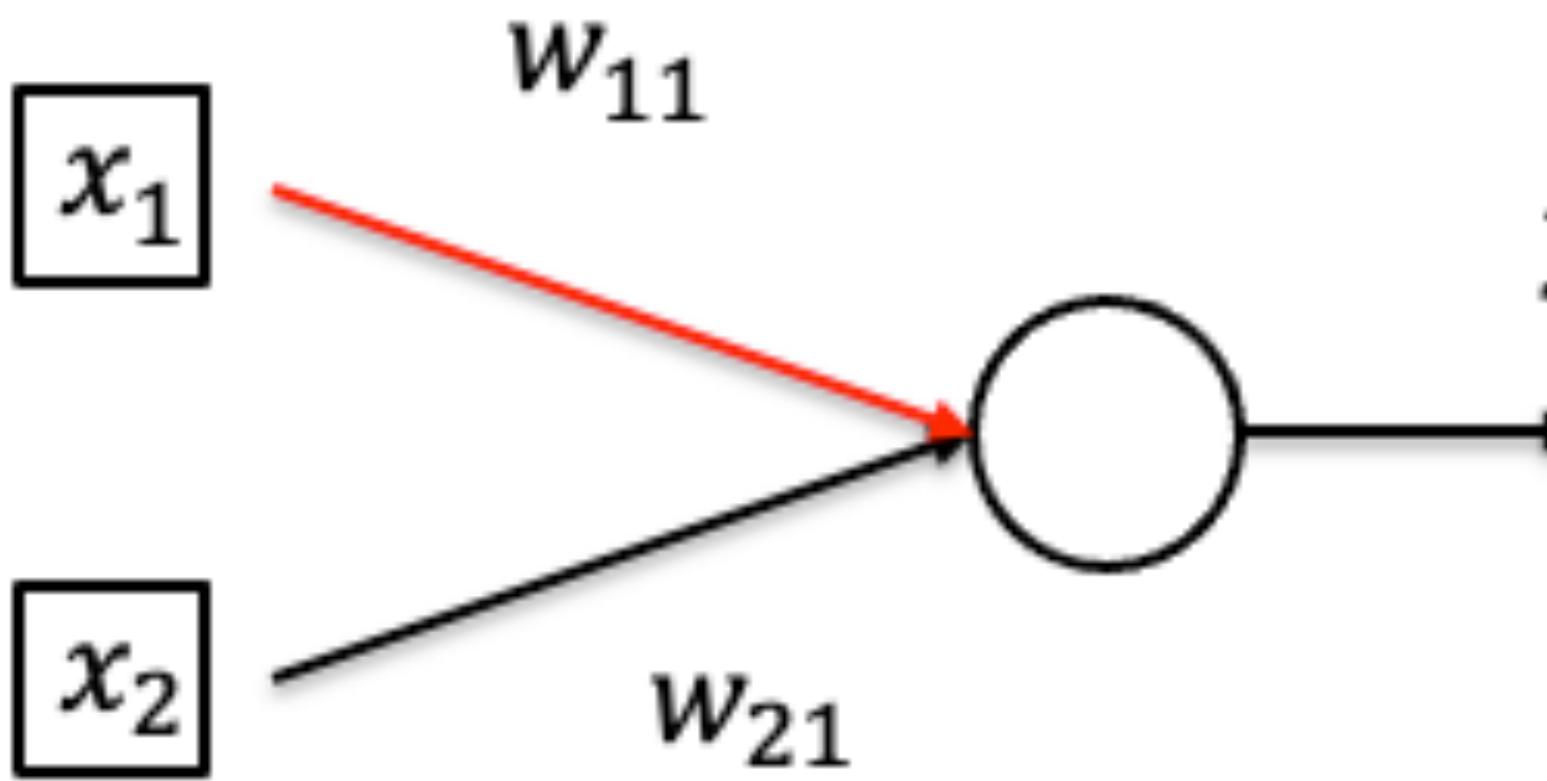
The backpropagation step shows the derivative of the loss with respect to the prediction \hat{y} :

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

- 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)

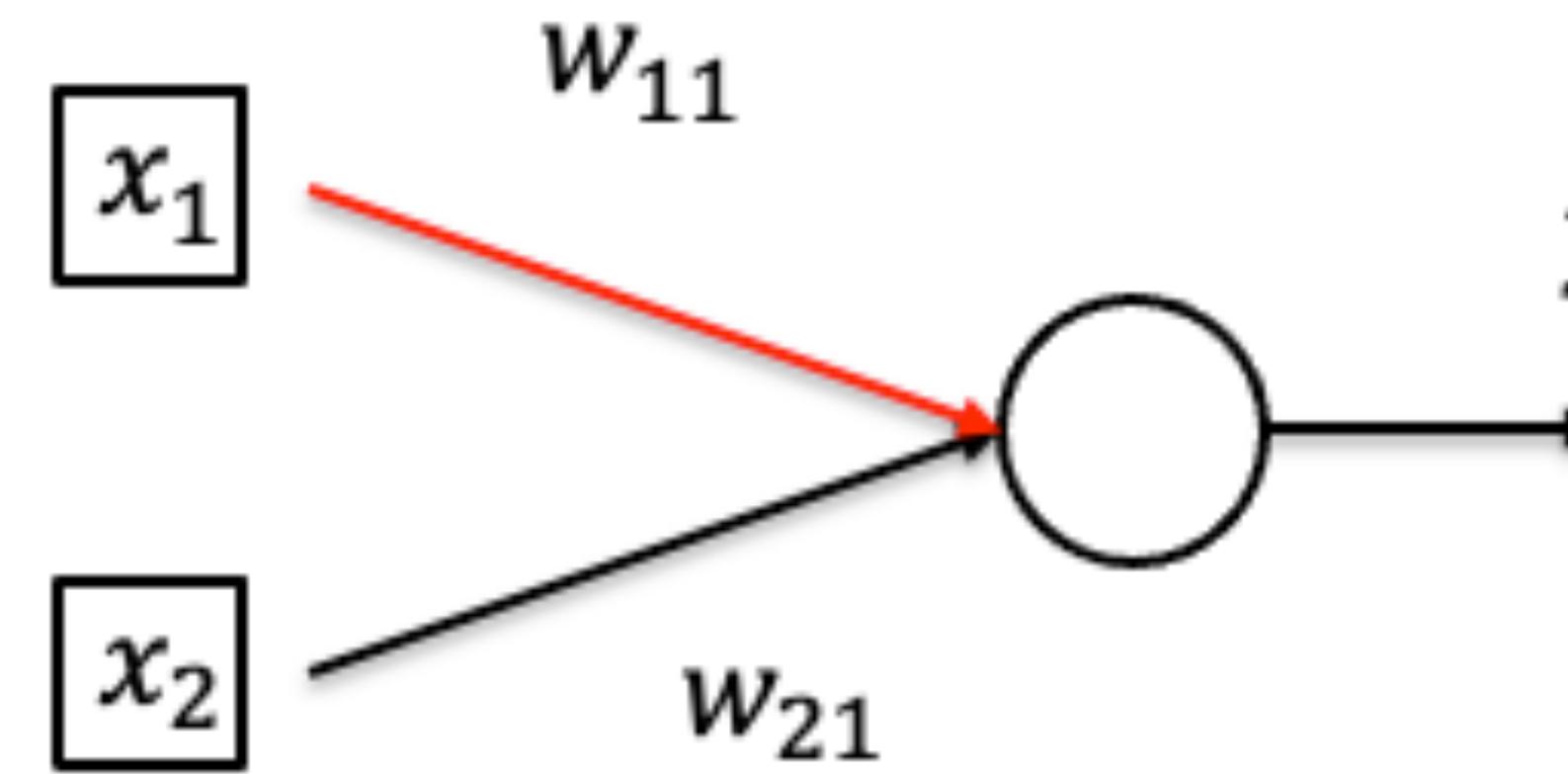


$$\begin{array}{ccccc} w_{11}x_1 & \xrightarrow{\text{+}} & z & \xrightarrow{\text{sigmoid function}} & \hat{y} \\ w_{21}x_2 & & & & \xrightarrow{-y \log(\hat{y})} \\ & & & & -(1 - \hat{y}) \log(1 - \hat{y}) \\ & & & & \ell(\mathbf{x}, y) \\ \frac{\partial \hat{y}}{\partial z} = \sigma'(z) & & & \frac{\partial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

- 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)

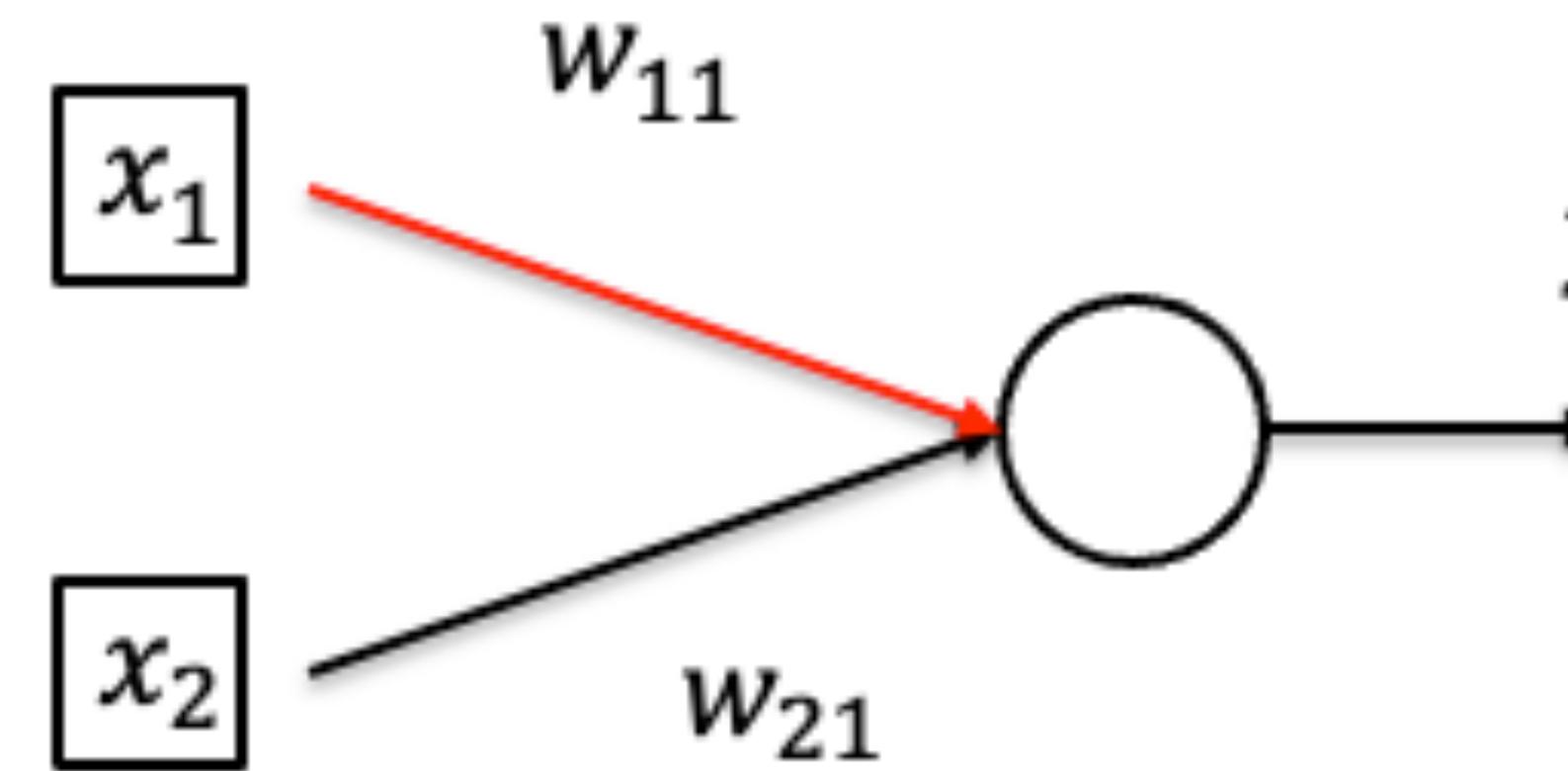


A computational graph illustrating the forward pass and the backward pass for calculating the gradient. The forward pass consists of three main steps: 1) Inputs $w_{11}x_1$ and $w_{21}x_2$ are summed at a plus node to produce z . 2) The sigmoid function maps z to \hat{y} . 3) The loss function maps \hat{y} to the loss value $\ell(\mathbf{x}, y)$. The backward pass involves calculating gradients: $\frac{\partial \hat{y}}{\partial z} = \sigma'(z) = \sigma(z)(1 - \sigma(z))$, $\frac{\partial \ell}{\partial \hat{y}} = -\frac{y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})}{\hat{y}(1 - \hat{y})}$, and $\frac{\partial l}{\partial w_{11}} = \frac{\partial l}{\partial \hat{y}} \cdot \hat{y}(1 - \hat{y})x_1$.

- By chain rule:

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

Calculate Gradient (on one data point)



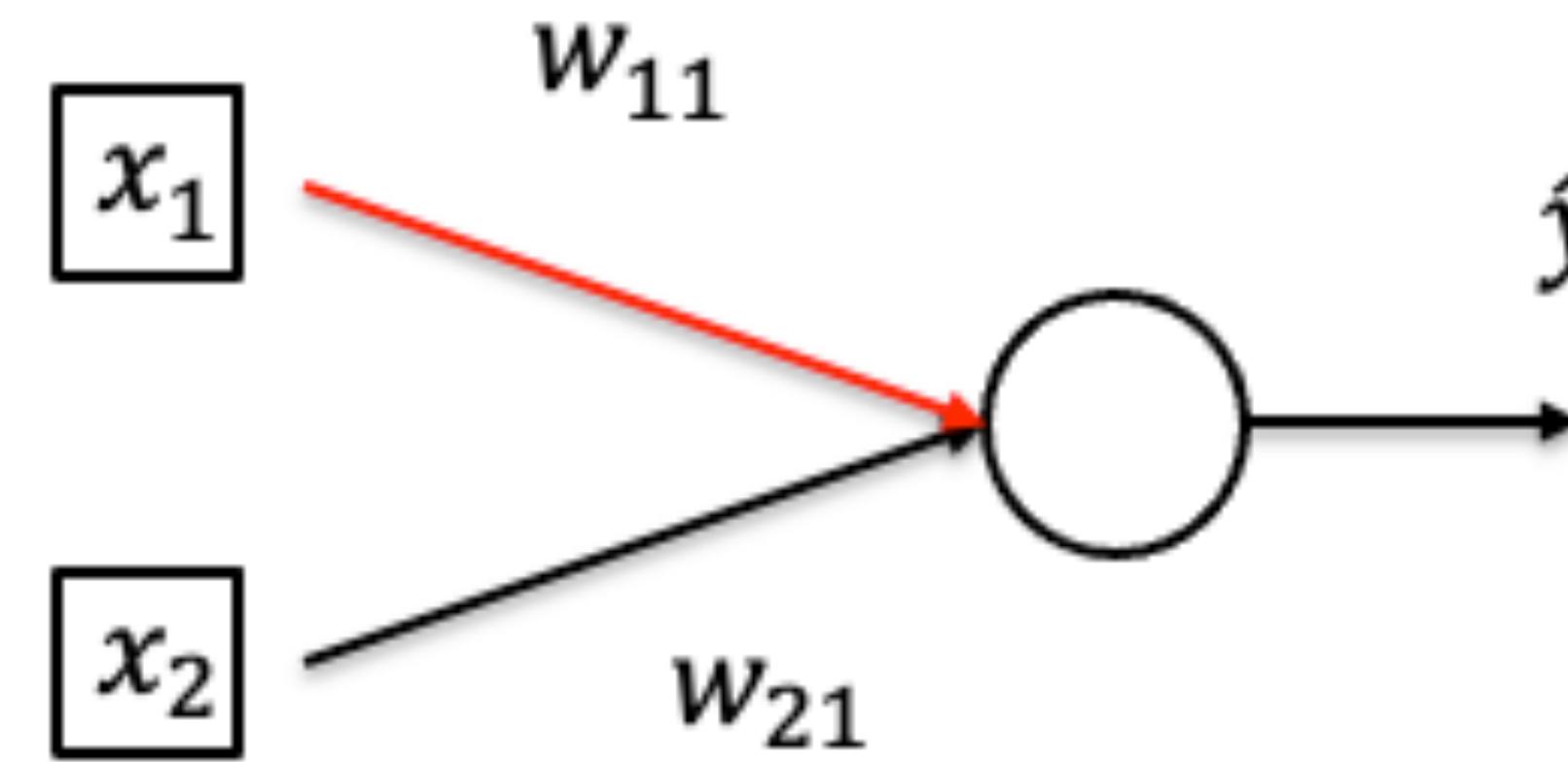
The diagram illustrates the forward pass of a neural network. It starts with inputs $w_{11}x_1$ and $w_{21}x_2$ which are summed at a plus sign. The result is passed through a "sigmoid function" to produce the output \hat{y} . The output \hat{y} is then used in the loss function $\ell(\mathbf{x}, y)$. Below the diagram, the derivative of the sigmoid function is given as $\frac{\partial \hat{y}}{\partial z} = \sigma'(z) = \sigma(z)(1 - \sigma(z))$.

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

- 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)



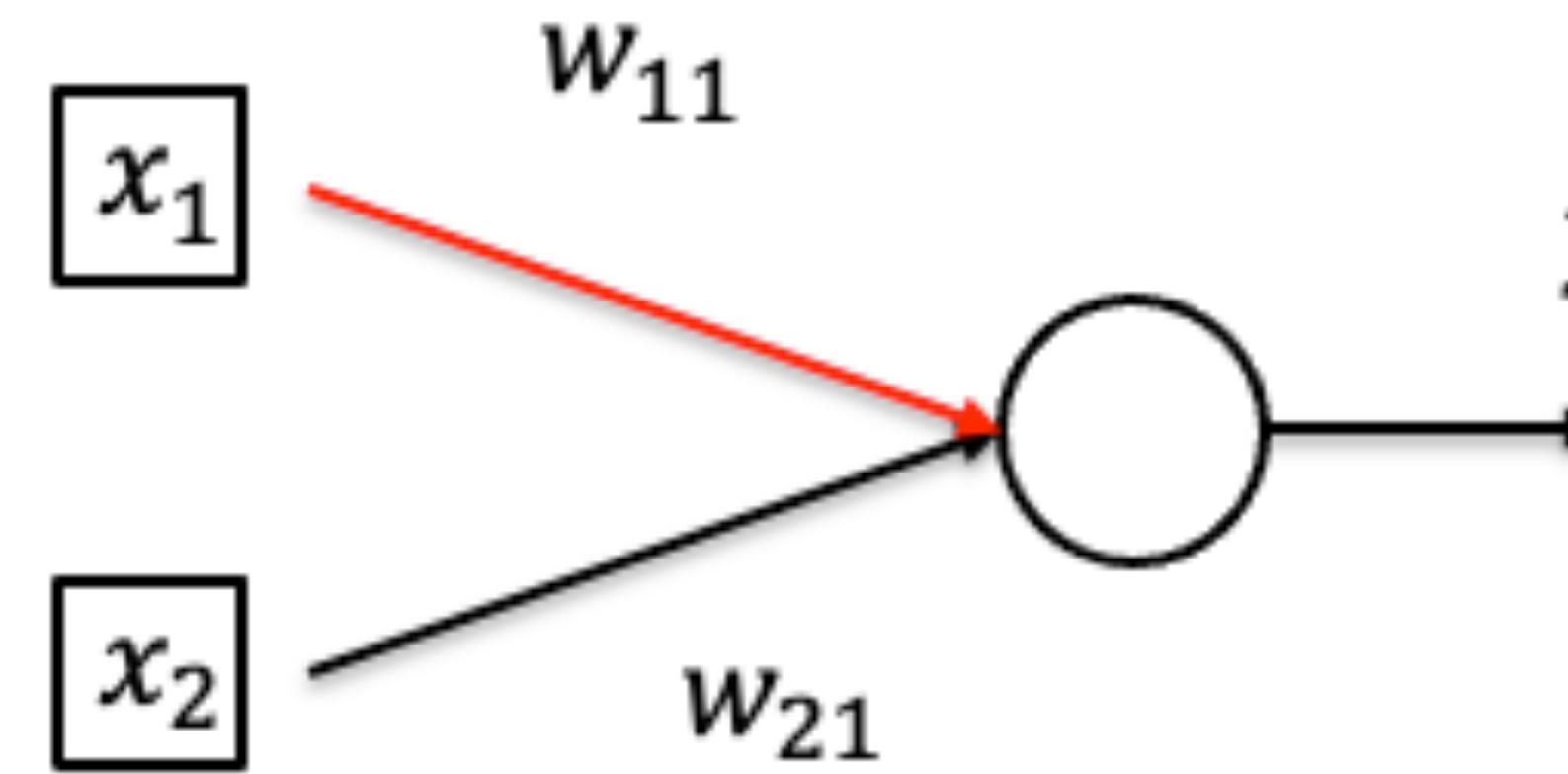
A computational graph illustrating the forward pass and the calculation of gradients. On the left, two inputs $w_{11}x_1$ and $w_{21}x_2$ are summed to produce z . This z is passed through a "sigmoid function" to produce the output \hat{y} . The loss function is given as $\ell(\mathbf{x}, y) = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y})$. Below this, the derivative of the sigmoid function is shown: $\frac{\partial \hat{y}}{\partial z} = \sigma'(z) = \sigma(z)(1 - \sigma(z))$.

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

- By chain rule:

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

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 .

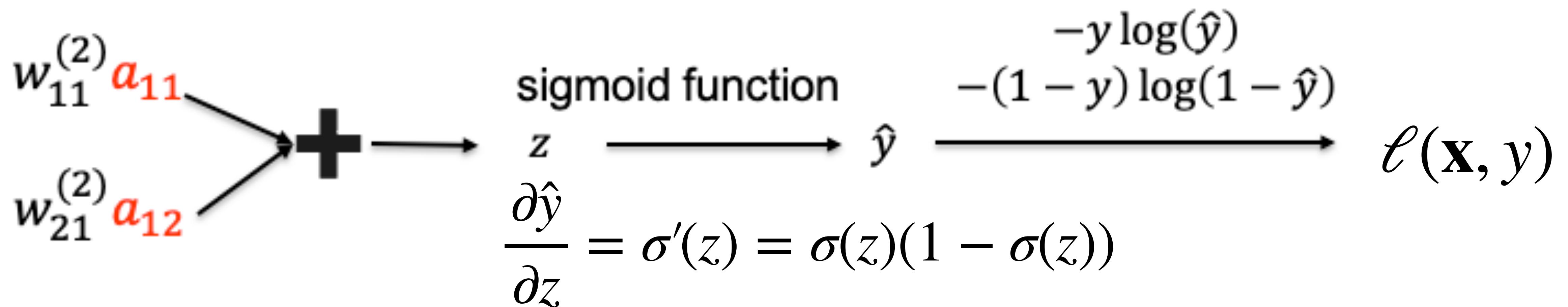
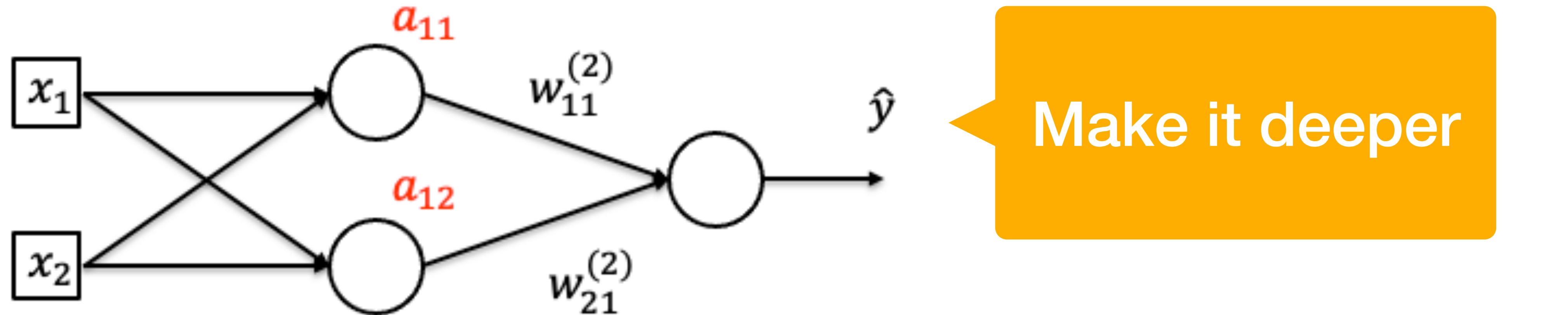
Below the graph, the gradient of the loss function with respect to the input x_1 is derived using the chain rule:

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

- By chain rule:

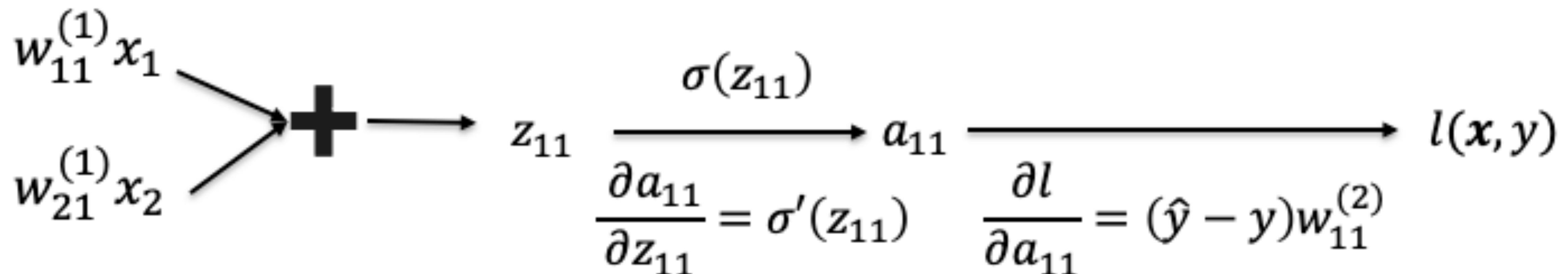
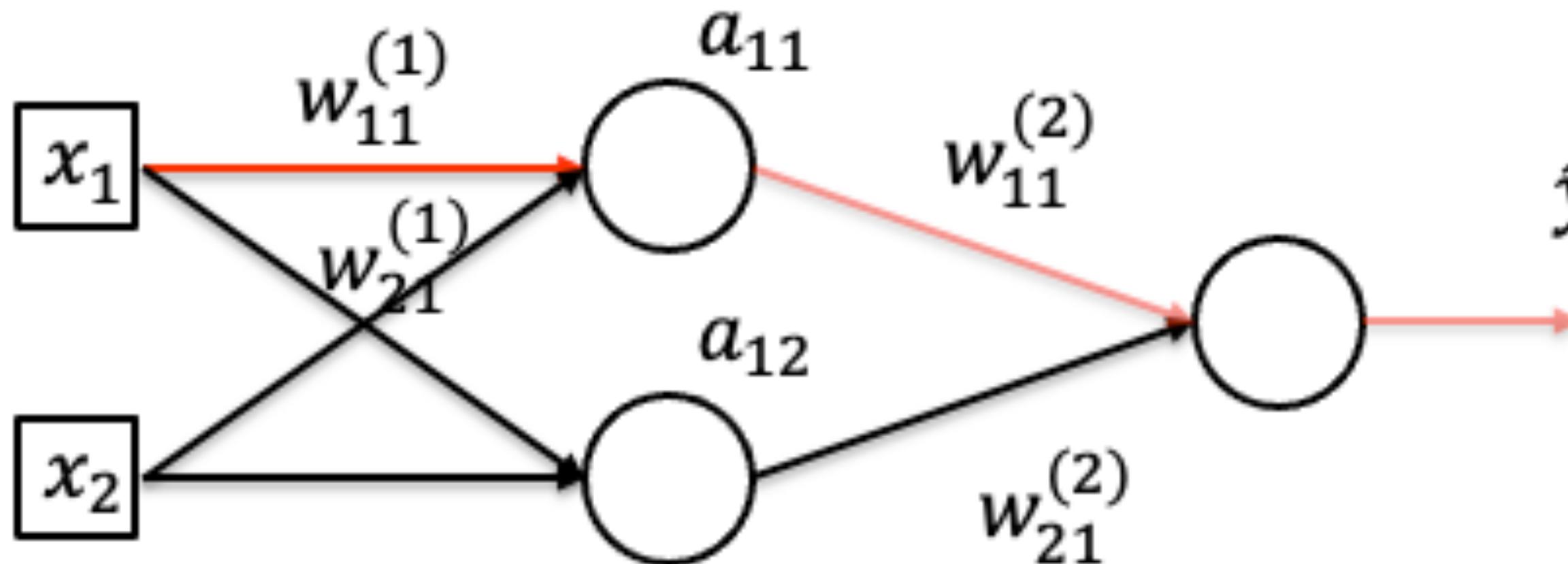
$$\frac{\partial l}{\partial x_1} = \frac{\partial l}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} w_{11} = (\hat{y} - y)w_{11}$$

Calculate Gradient (on one data point)



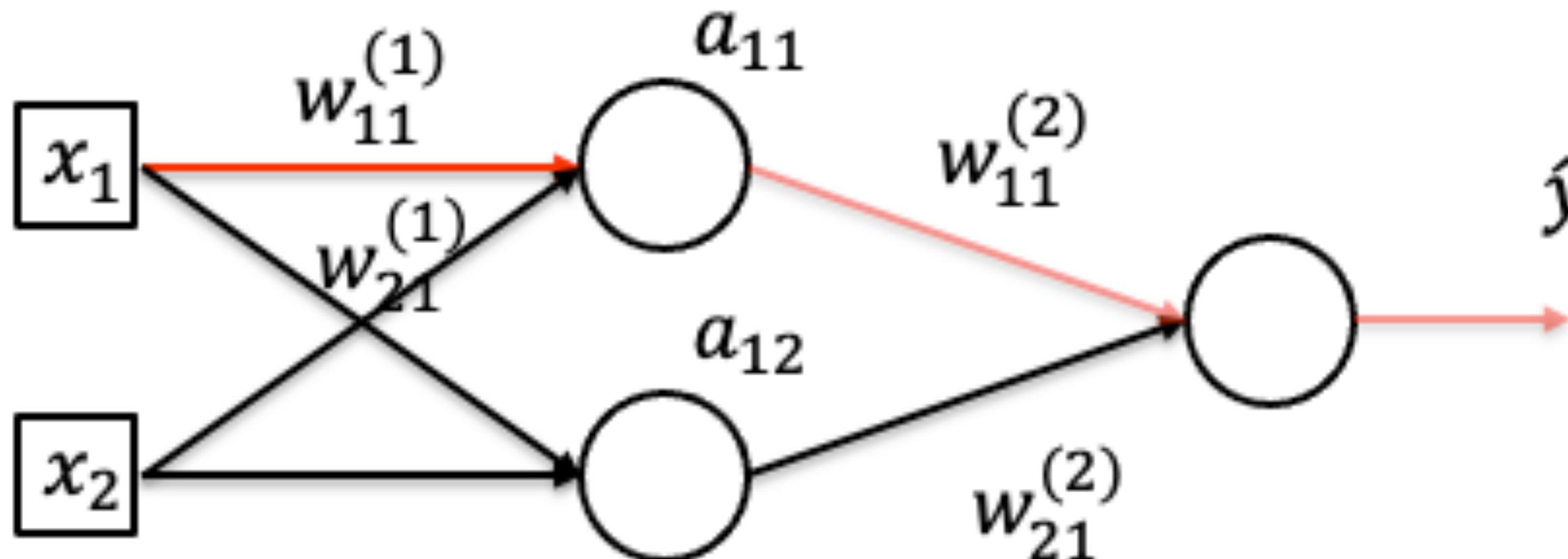
- By chain rule: $\frac{\partial l}{\partial a_{11}} = (\hat{y} - y)w_{11}^{(2)}$, $\frac{\partial l}{\partial a_{12}} = (\hat{y} - y)w_{21}^{(2)}$

Calculate Gradient (on one data point)



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

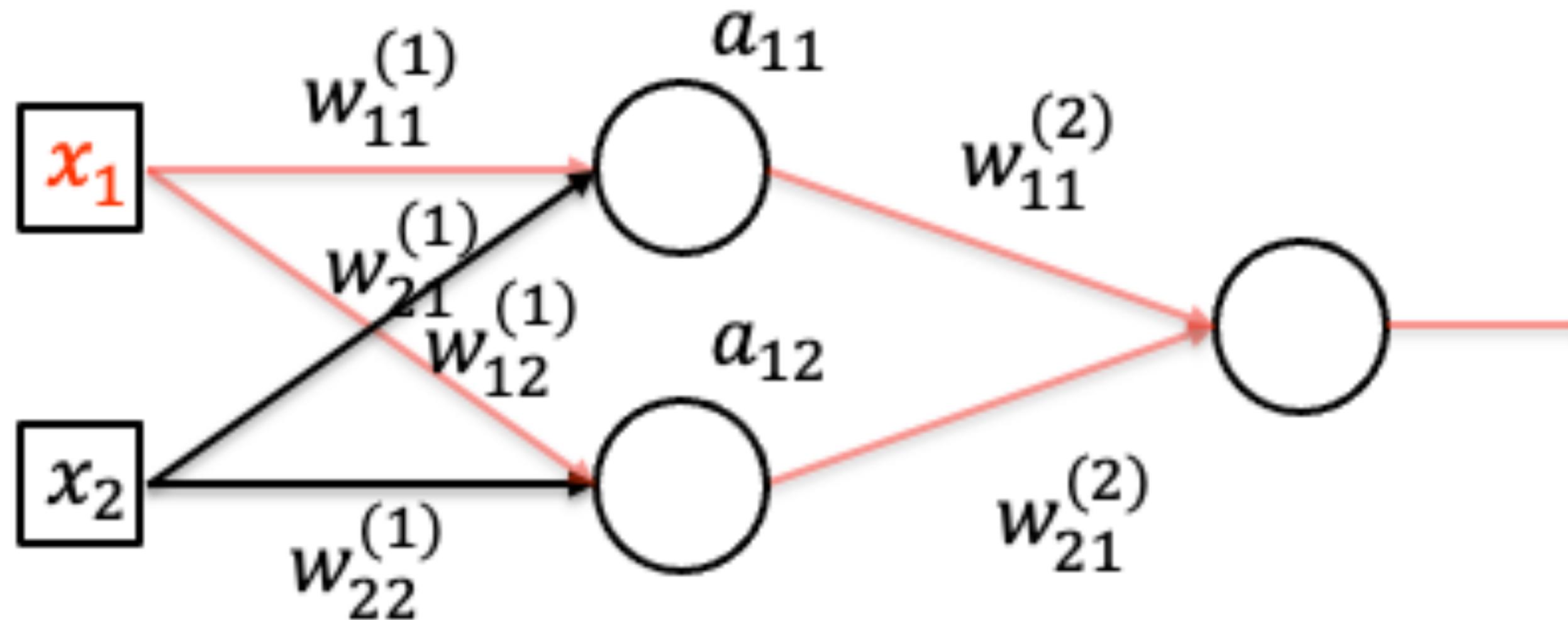
Calculate Gradient (on one data point)



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

- 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)



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

- 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

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

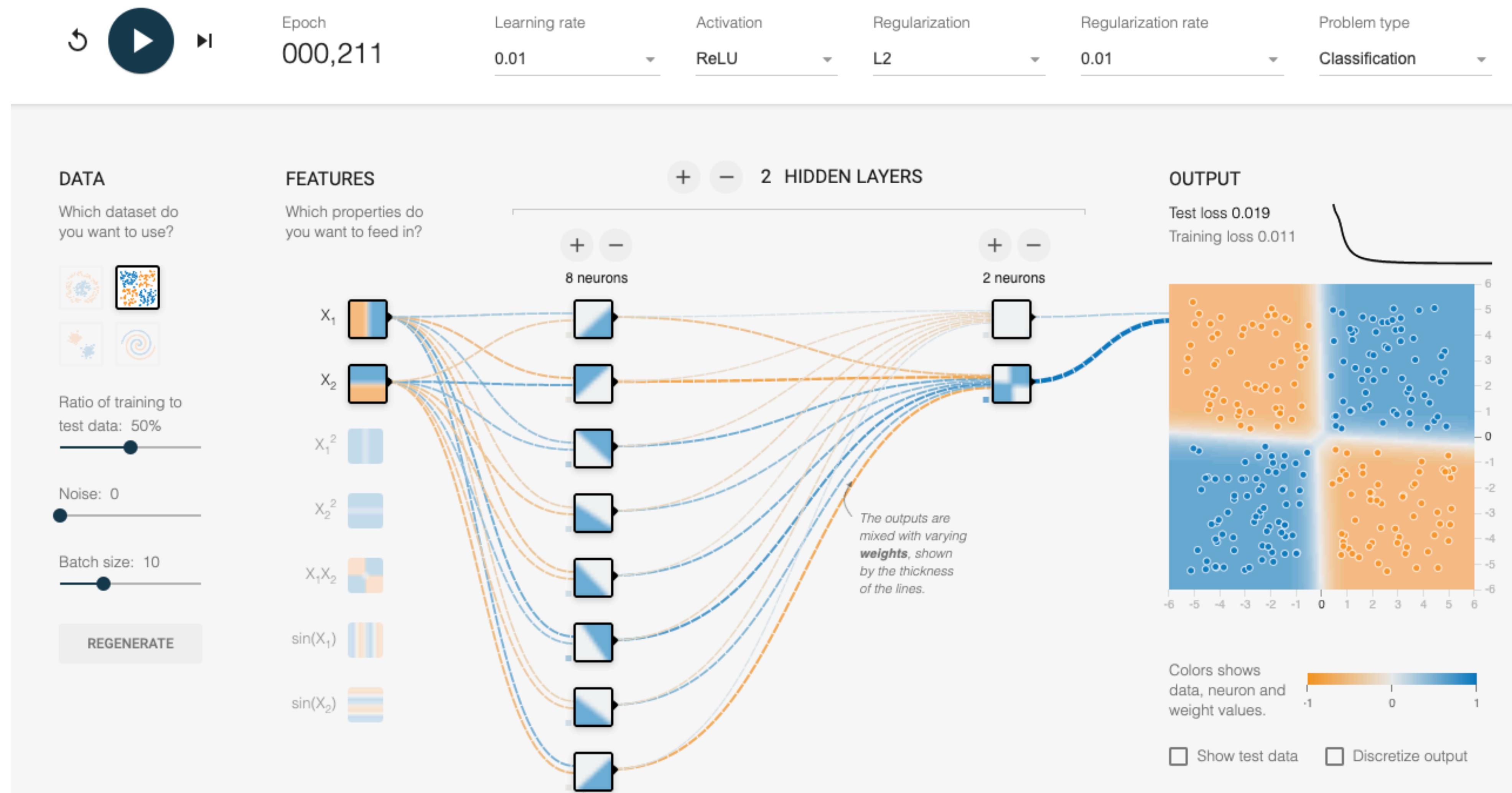
- A gradients, parameters
- B parameters, gradients
- C loss, parameters
- D parameters, loss

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 enough accuracy is needed?

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

Demo: Learning XOR using neural net



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

What we've learned today...

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