

Neural Networks

Part 2

Yin Li

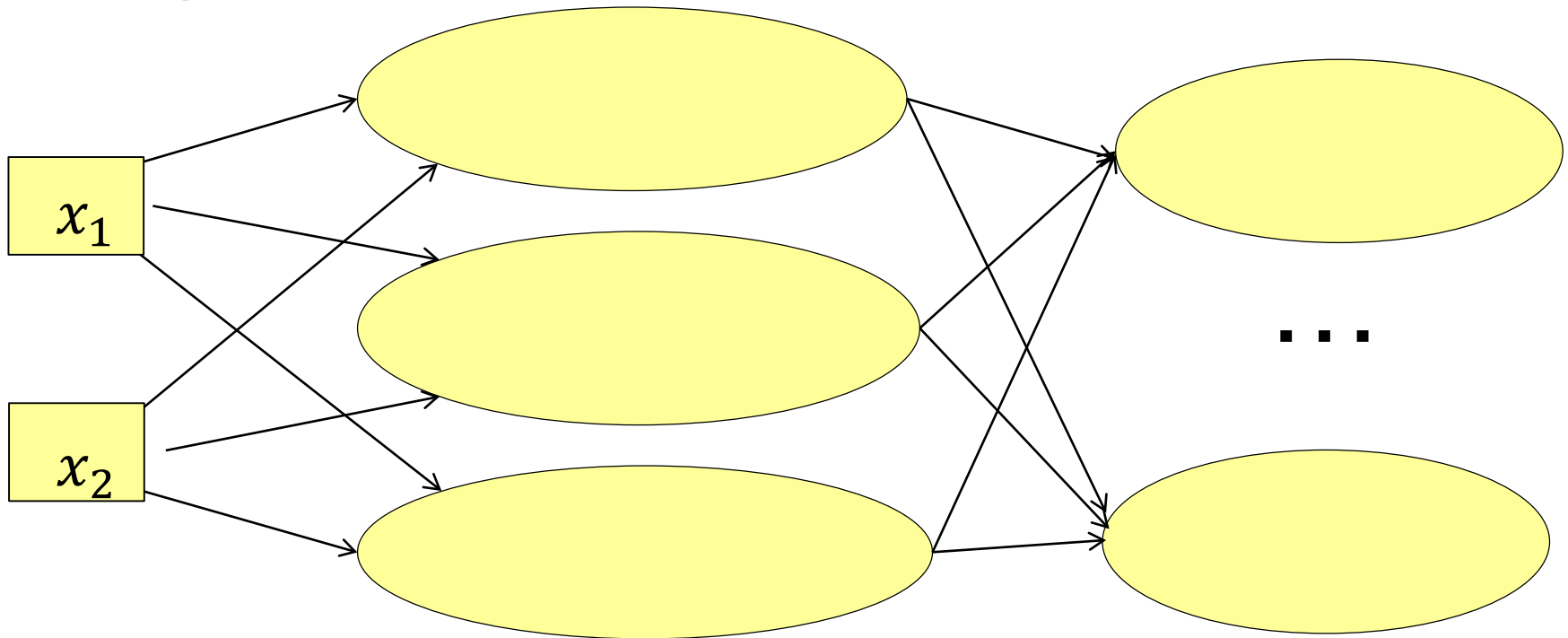
`yin.li@wisc.edu`

University of Wisconsin, Madison

[Some of the slides from Yingyu Liang, Jerry Zhu, Mohit Gupta]

Neural net for K -way classification

- Use K output units
- Training: encode a label y by an indicator vector
 - class1=(1,0,0,...,0), class2=(0,1,0,...,0) etc.
- Test: choose the class corresponding to the largest output unit



Mini-batch stochastic gradient descent

- Define a loss function $E_x = E(x, y)$
- Select a learning rate $\alpha > 0$
- Initialize the model parameters (edge weights) $w^{(0)}$
- For $t = 1, 2, \dots$
 - Randomly sample a subset \hat{D} from D
 - Compute $\frac{\partial E_x}{\partial w}$ (per sample gradients w.r.t. w) for $x \in \hat{D}$
 - Update the parameters

$$w^{(t)} = w^{(t-1)} - \alpha \frac{1}{|\hat{D}|} \sum_{x \in \hat{D}} \frac{\partial E_x}{\partial w}$$

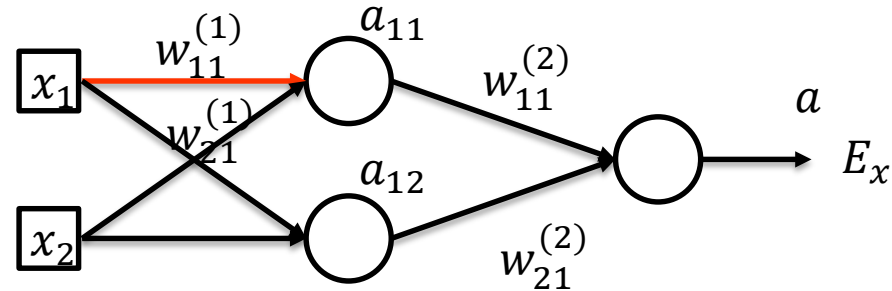
- Repeat until E converges

The **key challenge** is to compute $\frac{\partial E_x}{\partial w}$!

Back-propagation

- In theory

$$\frac{\partial E_x}{\partial w_{11}} = (a - y)w_{11}^{(2)} a_{11} (1 - a_{11})x_1$$

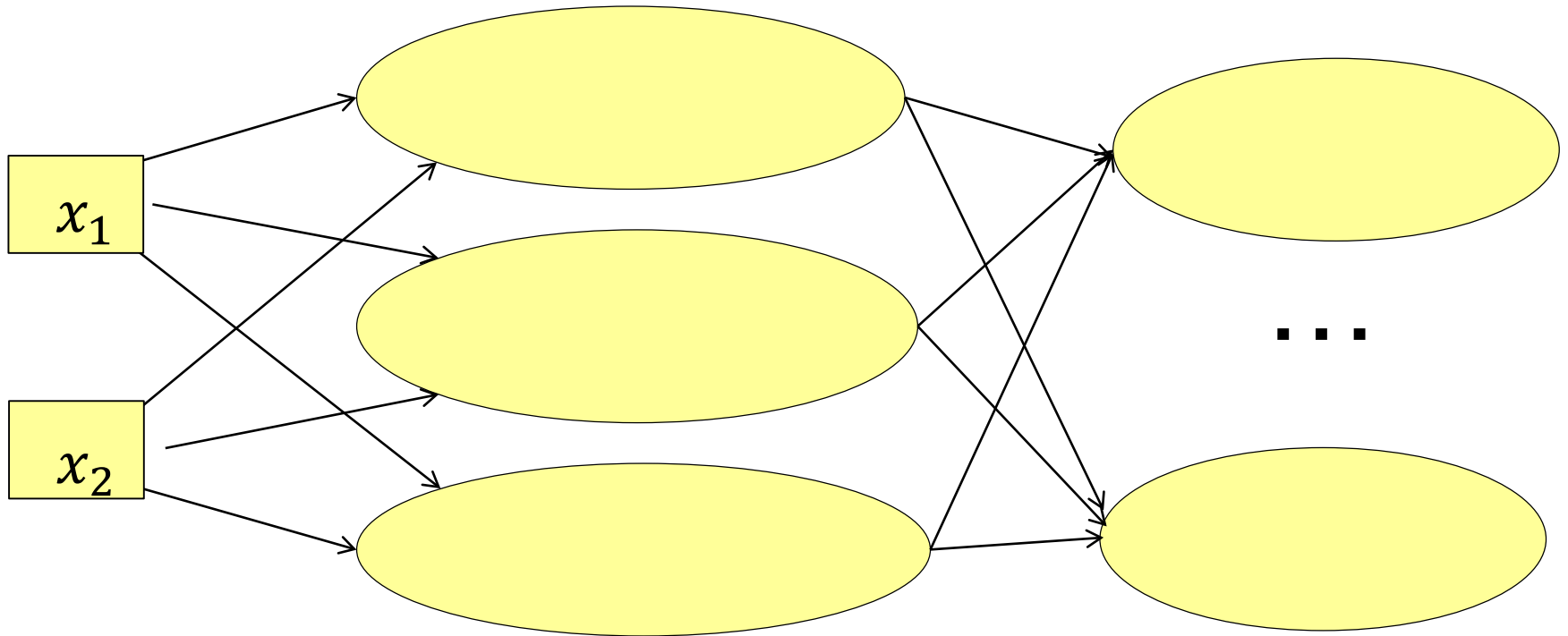


- In practice

- Define the model and the loss function
- Select the optimization method (e.g., stochastic gradient descent)
- Back-propagation handled by **automatic differentiation**
- How?

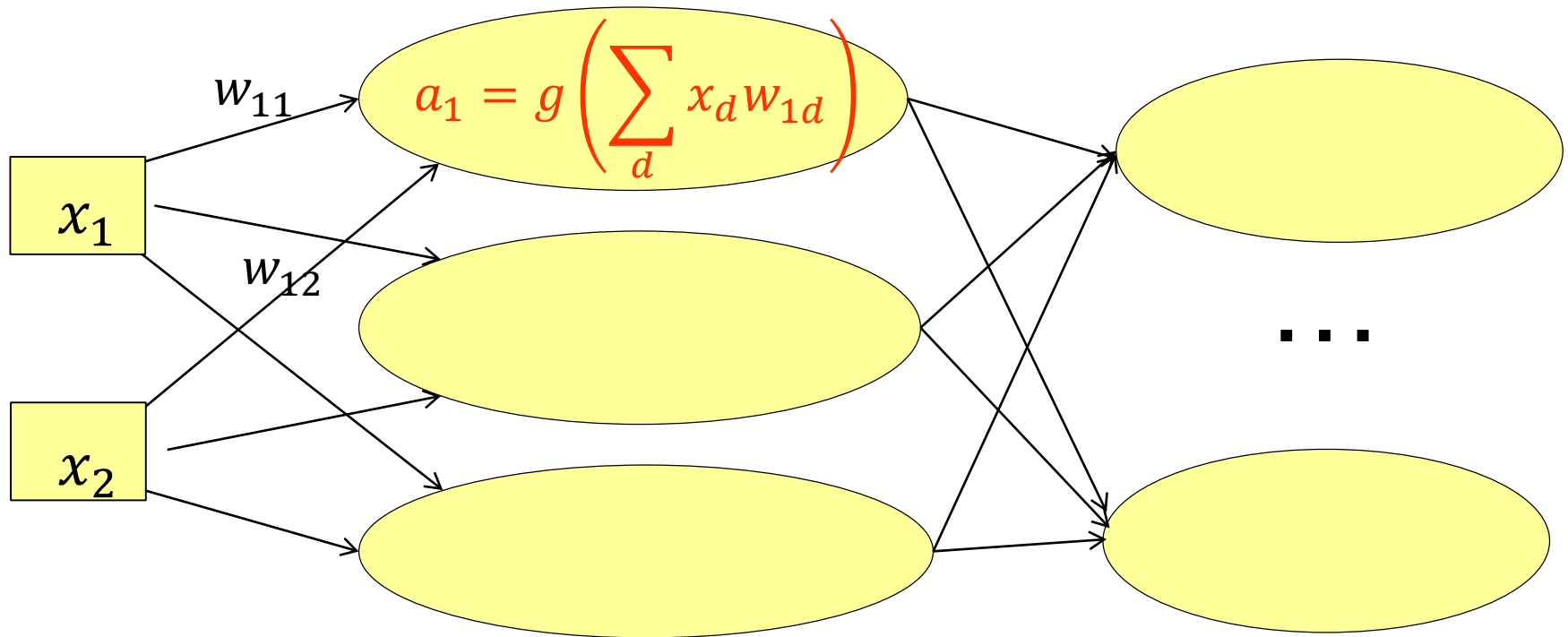
A single layer in neural network

- Let $\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$



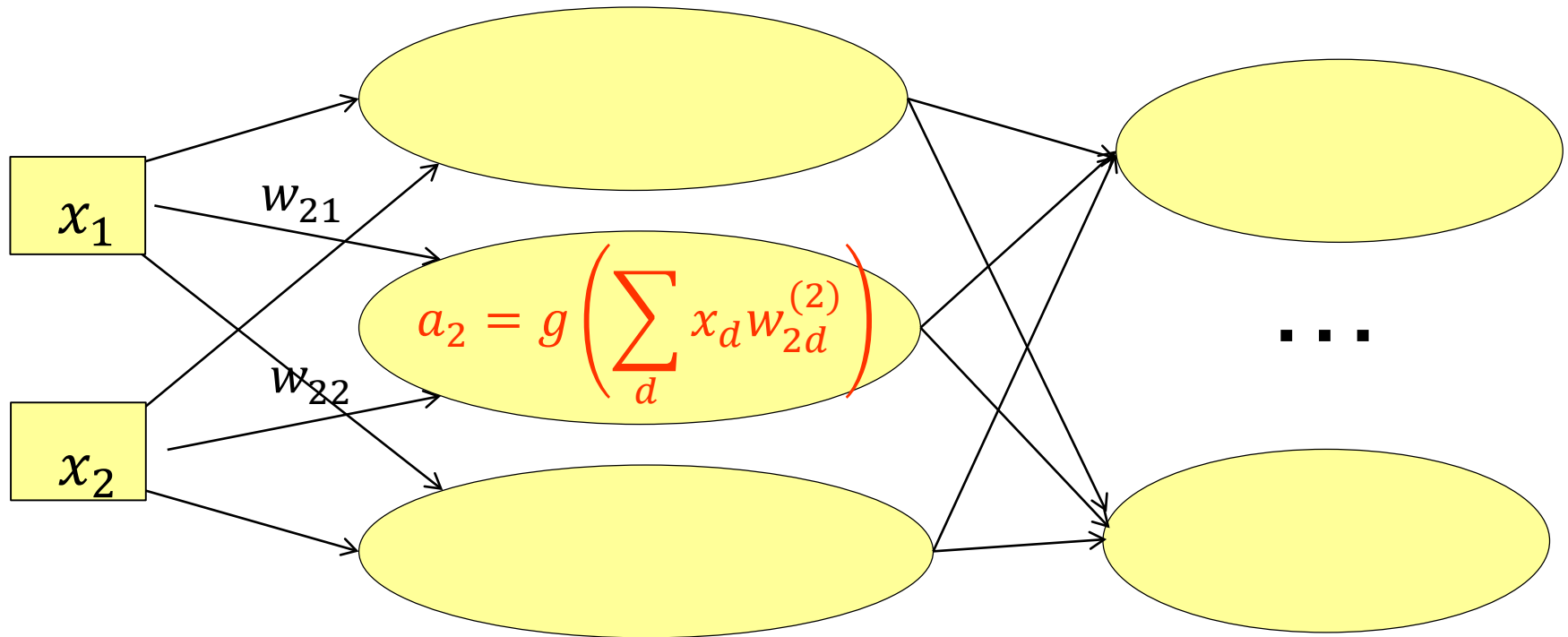
A single layer in neural network

- Let $\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$
- $a_1 = g \left([w_{11} \ w_{12}] * \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + w_{10} \right) = g(\mathbf{w}_1^T \mathbf{x} + b_1)$



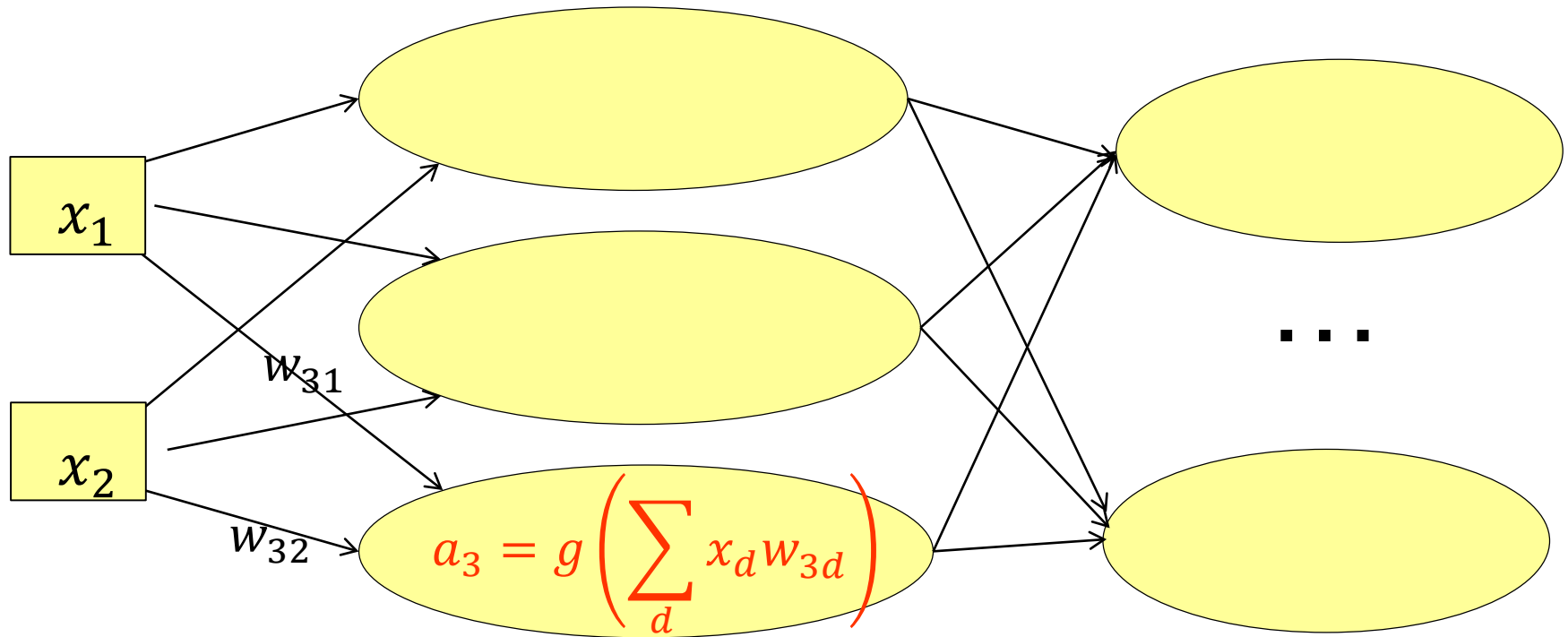
A single layer in neural network

- Let $\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$
- $a_2 = g \left([w_{21} \ w_{22}] * \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + w_{20} \right) = g(\mathbf{w}_2^T \mathbf{x} + b_2)$



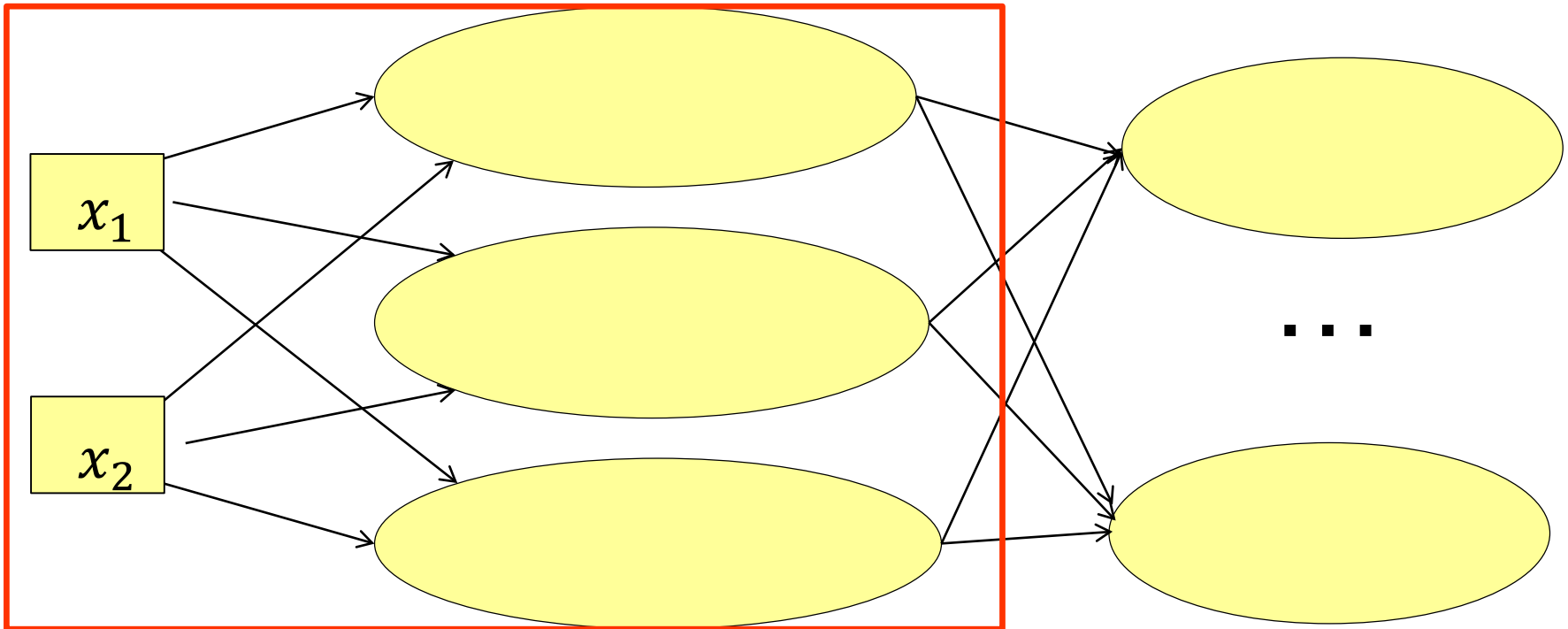
A single layer in neural network

- Let $\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$
- $a_3 = g \left([w_{31} \ w_{32}] * \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + w_{30} \right) = g(\mathbf{w}_3^T \mathbf{x} + b_3)$



A single layer in neural network

- $a_1 = g\left([w_{11} \ w_{12}] * \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + w_{10}\right) = g(\mathbf{w}_1^T \mathbf{x} + b_1)$
- $a_2 = g\left([w_{21} \ w_{22}] * \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + w_{20}\right) = g(\mathbf{w}_2^T \mathbf{x} + b_2)$
- $a_3 = g\left([w_{31} \ w_{32}] * \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + w_{30}\right) = g(\mathbf{w}_3^T \mathbf{x} + b_3)$



A single layer in neural network

- $a_1 = g\left([w_{11} \ w_{12}] * \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + w_{10}\right) = g(\mathbf{w}_1^T \mathbf{x} + b_1)$
- $a_2 = g\left([w_{21} \ w_{22}] * \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + w_{20}\right) = g(\mathbf{w}_2^T \mathbf{x} + b_2)$
- $a_3 = g\left([w_{31} \ w_{32}] * \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + w_{30}\right) = g(\mathbf{w}_3^T \mathbf{x} + b_3)$

g : Any activation function that applies on every element of the input (element-wise operation)

$$\mathbf{a} = \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix} = \begin{bmatrix} g(\mathbf{w}_1^T \mathbf{x} + b_1) \\ g(\mathbf{w}_2^T \mathbf{x} + b_2) \\ g(\mathbf{w}_3^T \mathbf{x} + b_3) \end{bmatrix} = g\left(\begin{bmatrix} \mathbf{w}_1^T \mathbf{x} + b_1 \\ \mathbf{w}_2^T \mathbf{x} + b_2 \\ \mathbf{w}_3^T \mathbf{x} + b_3 \end{bmatrix}\right)$$

A single layer in neural network

- $a_1 = g\left([w_{11} \ w_{12}] * \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + w_{10}\right) = g(\mathbf{w}_1^T \mathbf{x} + b_1)$
- $a_2 = g\left([w_{21} \ w_{22}] * \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + w_{20}\right) = g(\mathbf{w}_2^T \mathbf{x} + b_2)$
- $a_3 = g\left([w_{31} \ w_{32}] * \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + w_{30}\right) = g(\mathbf{w}_3^T \mathbf{x} + b_3)$

Rearrange the linear operations

$$\mathbf{a} = \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix} = g\left(\begin{bmatrix} \mathbf{w}_1^T \mathbf{x} + b_1 \\ \mathbf{w}_2^T \mathbf{x} + b_2 \\ \mathbf{w}_3^T \mathbf{x} + b_3 \end{bmatrix}\right) = g\left(\begin{bmatrix} \mathbf{w}_1^T \\ \mathbf{w}_2^T \\ \mathbf{w}_3^T \end{bmatrix} \mathbf{x} + \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}\right)$$

A single layer in neural network

- $a_1 = g\left([w_{11} \ w_{12}] * \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + w_{10}\right) = g(\mathbf{w}_1^T \mathbf{x} + b_1)$
- $a_2 = g\left([w_{21} \ w_{22}] * \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + w_{20}\right) = g(\mathbf{w}_2^T \mathbf{x} + b_2)$
- $a_3 = g\left([w_{31} \ w_{32}] * \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + w_{30}\right) = g(\mathbf{w}_3^T \mathbf{x} + b_3)$

Rearrange the linear operations

$$\mathbf{a} = \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix} = g\left(\begin{bmatrix} \mathbf{w}_1^T \\ \mathbf{w}_2^T \\ \mathbf{w}_3^T \end{bmatrix} \mathbf{x} + \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}\right) = g(\mathbf{W}^T \mathbf{x} + \mathbf{b})$$

A single layer in neural network

- $a_1 = g\left([w_{11} \ w_{12}] * \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + w_{10}\right) = g(\mathbf{w}_1^T \mathbf{x} + b_1)$
- $a_2 = g\left([w_{21} \ w_{22}] * \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + w_{20}\right) = g(\mathbf{w}_2^T \mathbf{x} + b_2)$
- $a_3 = g\left([w_{31} \ w_{32}] * \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + w_{30}\right) = g(\mathbf{w}_3^T \mathbf{x} + b_3)$

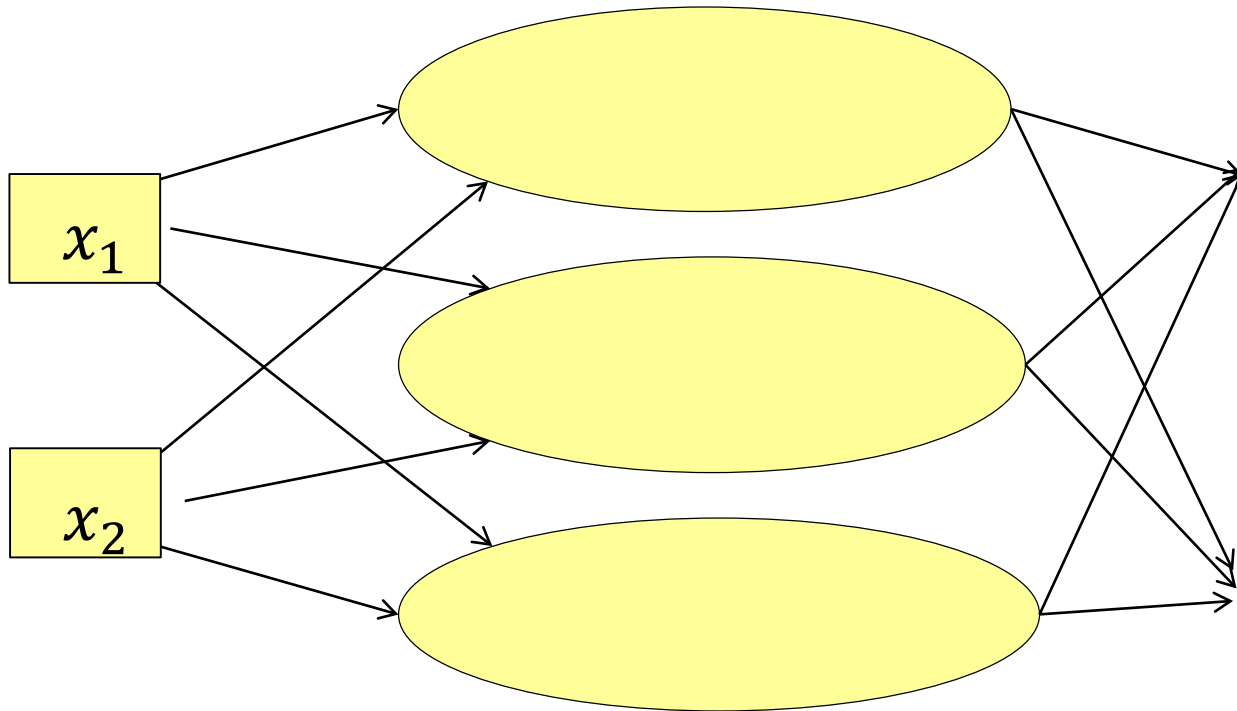
Here is what we have

$$\mathbf{a} = g(\mathbf{W}^T \mathbf{x} + \mathbf{b})$$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad \mathbf{W} = \begin{bmatrix} w_{11} & w_{21} & w_{31} \\ w_{12} & w_{22} & w_{32} \end{bmatrix} \quad \mathbf{b} = \begin{bmatrix} w_{10} \\ w_{20} \\ w_{30} \end{bmatrix}$$

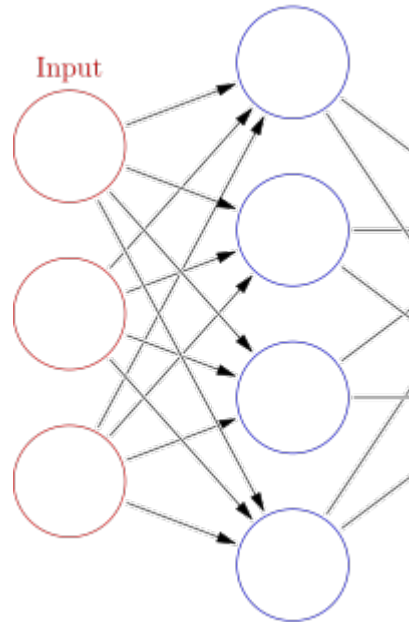
A single layer in neural network

- $a = g(W^T x + b)$
- Work for any element-wise activation function g



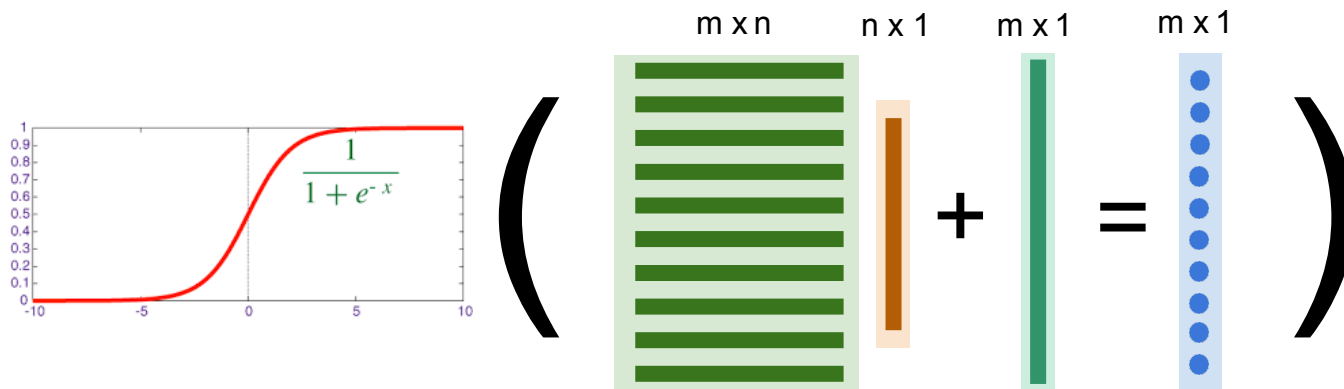
A single layer in neural network

- $\mathbf{a} = g(\mathbf{W}^T \mathbf{x} + \mathbf{b})$
- Work for any element-wise activation function g
- Work for any input / output dimensions
- Map an input $\mathbf{x} \in R^n$ to an output $\mathbf{a} \in R^m$
- $\mathbf{x} \in R^n$, $\mathbf{W} \in R^{n \times m}$, $\mathbf{b} \in R^m$, $\mathbf{a} \in R^m$



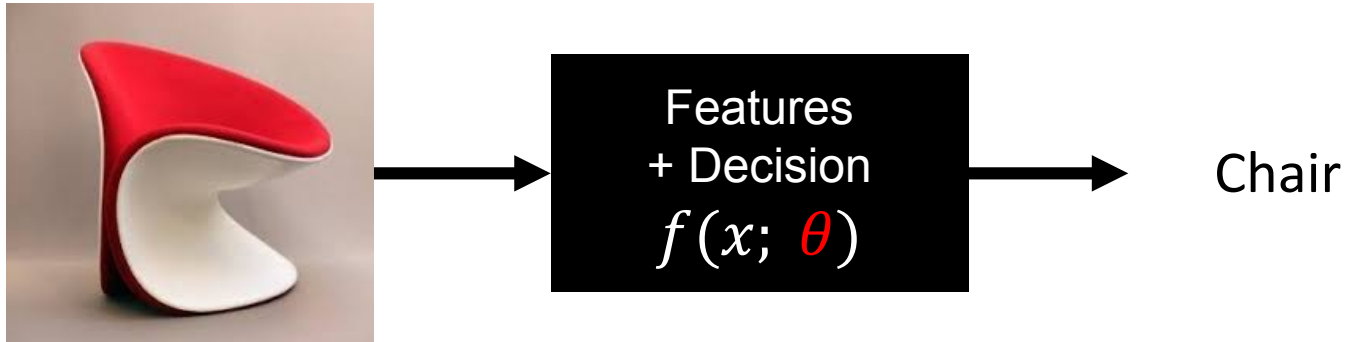
A single layer in neural network

- $\mathbf{a} = g(\mathbf{W}^T \mathbf{x} + \mathbf{b})$
- Work for any element-wise activation function g
- Work for any input / output dimensions
- Map an input $\mathbf{x} \in R^n$ to an output $\mathbf{a} \in R^m$
- $\mathbf{x} \in R^n$, $\mathbf{W} \in R^{n \times m}$, $\mathbf{b} \in R^m$, $\mathbf{a} \in R^m$



Neural Networks / Deep Learning

- What type of functions shall we consider for f ?



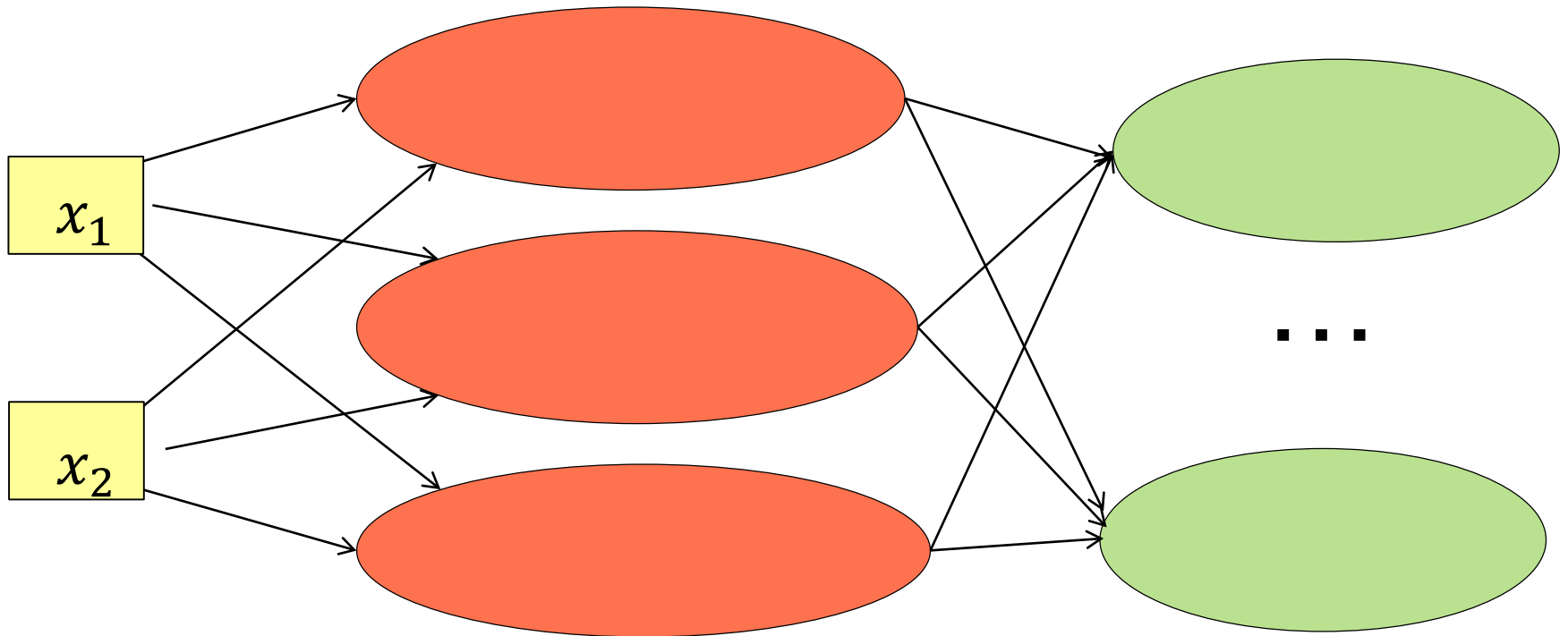
Proposal: Composing a set of (nonlinear) functions g

$$f(\mathbf{x}; \boldsymbol{\theta}) = g_1(\dots g_{n-1}(g_n(\mathbf{x}; \boldsymbol{\theta}_n), \boldsymbol{\theta}_{n-1}) \dots, \boldsymbol{\theta}_1)$$

Example: $a = \text{sigmoid}(\mathbf{W}^T \mathbf{x} + \mathbf{b}) = g(\mathbf{x}; \mathbf{W}, \mathbf{b})$

An example neural network

- Let $\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$
- $f(\mathbf{x}; \boldsymbol{\theta}) = g_1(g_2(\mathbf{x}; \boldsymbol{\theta}_2), \boldsymbol{\theta}_1)$



Neural Networks / Deep Learning

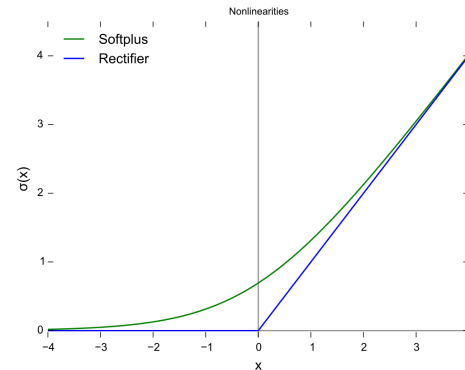
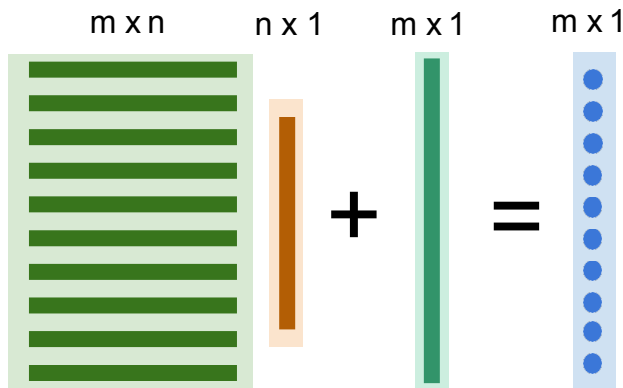
- (Deep) Neural Network:

Composing a set of (nonlinear) functions g

$$f(\mathbf{x}; \boldsymbol{\theta}) = g_1(\dots g_{n-1}(g_n(\mathbf{x}; \boldsymbol{\theta}_n), \boldsymbol{\theta}_{n-1}) \dots, \boldsymbol{\theta}_1)$$

- Key Elements:

Linear operations + Nonlinear activations



Neural Networks / Deep Learning

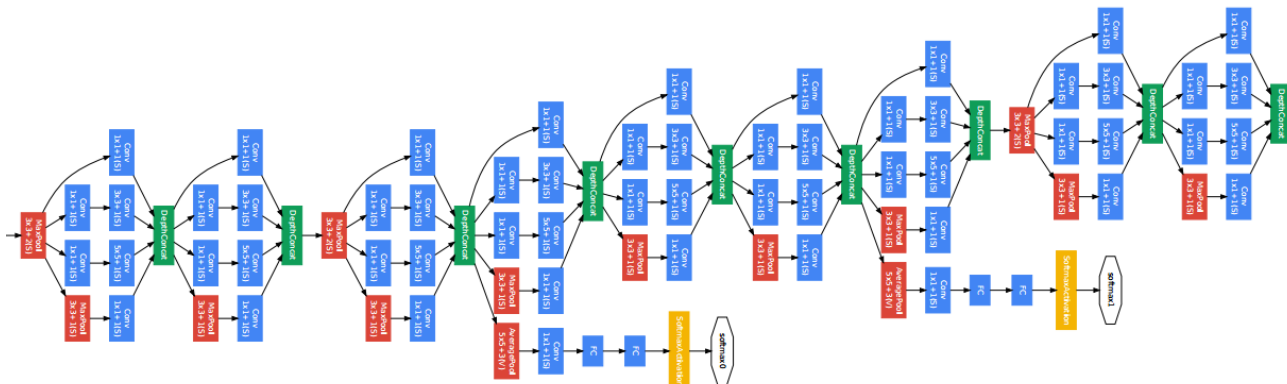
- (Deep) Neural Network:

Composing a set of (nonlinear) functions g

$$f(x; \theta) = g_1(\dots g_{n-1}(g_n(x; \theta_n), \theta_{n-1}) \dots, \theta_1)$$

- Key Challenge:

How to present the composition of nonlinear functions & make the computation possible and efficient?

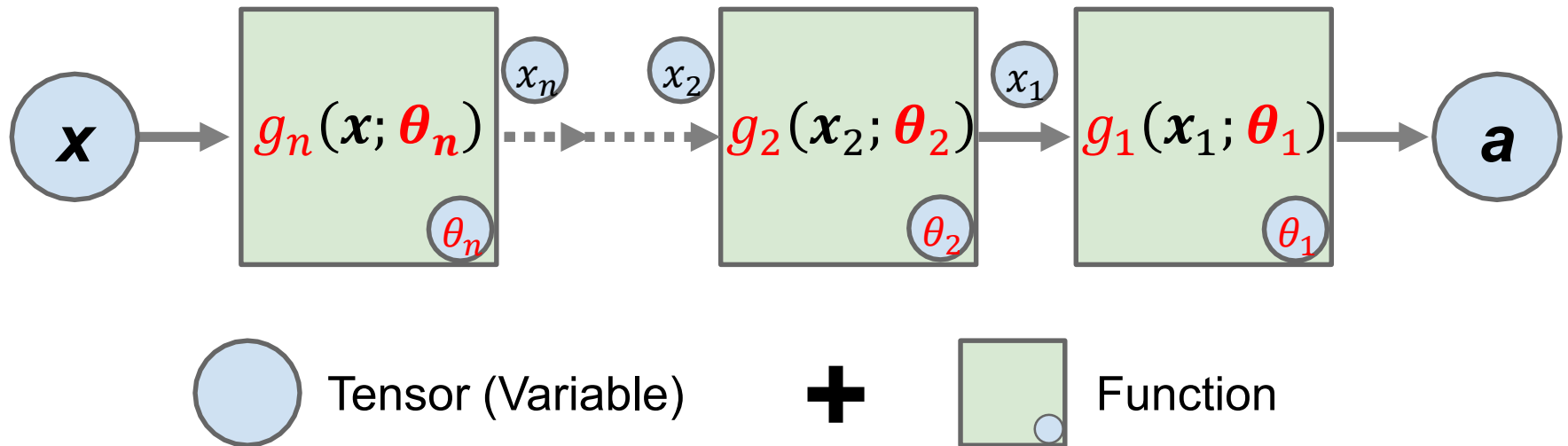


Neural network as a stack of layers

- (Deep) Neural Network:

Composing a set of (nonlinear) functions g

$$f(\mathbf{x}; \boldsymbol{\theta}) = g_1(\dots g_{n-1}(g_n(\mathbf{x}; \boldsymbol{\theta}_n), \boldsymbol{\theta}_{n-1}) \dots, \boldsymbol{\theta}_1)$$

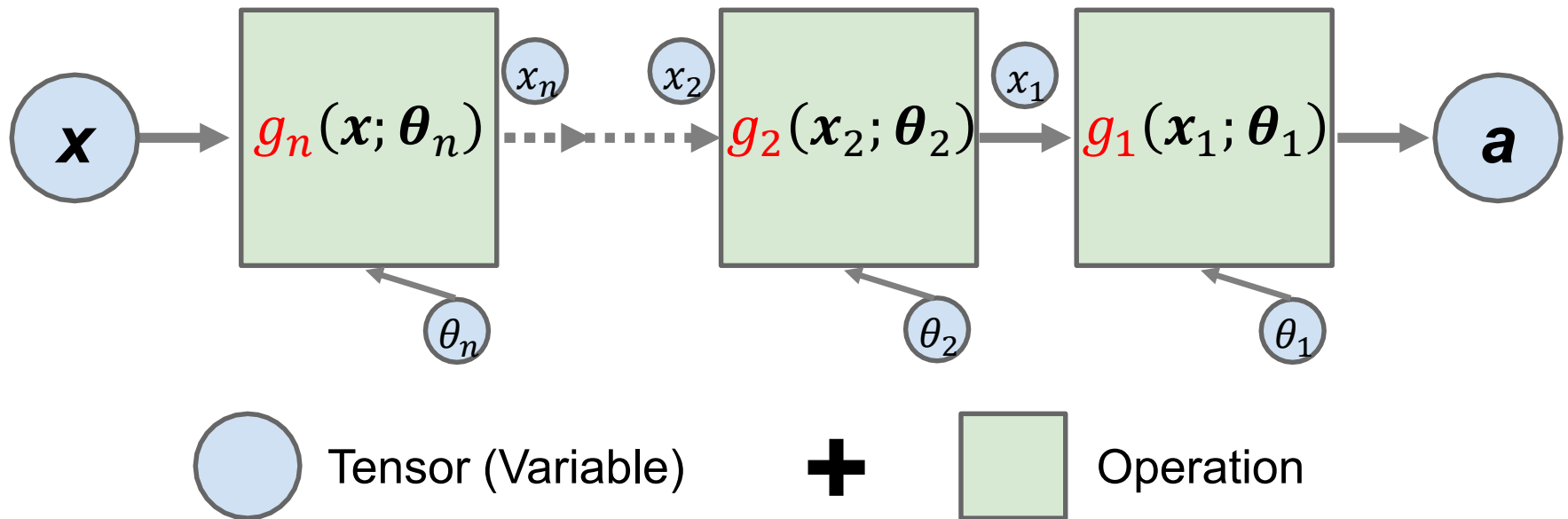


Neural network as variables + operations

- (Deep) Neural Network:

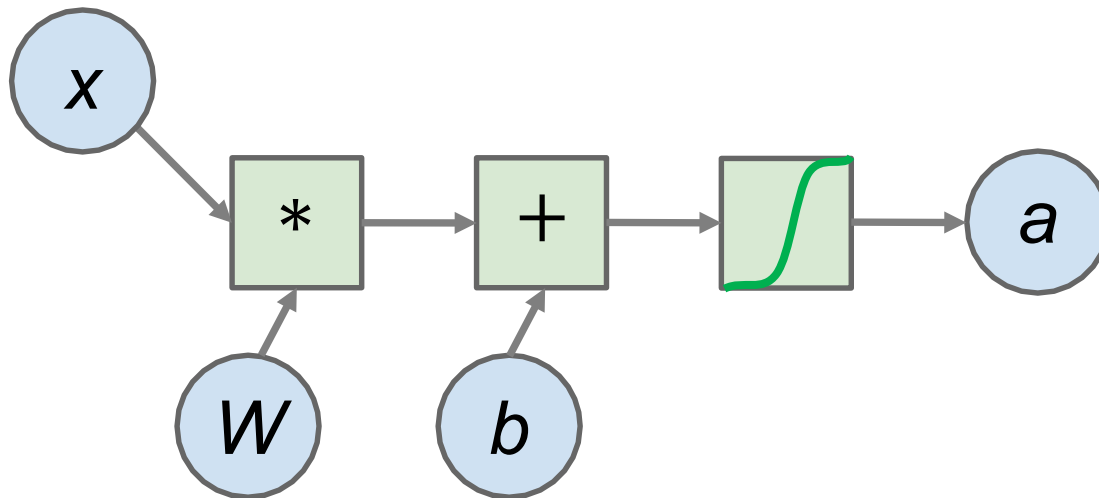
Composing a set of (nonlinear) functions g

$$f(\mathbf{x}; \boldsymbol{\theta}) = g_1(\dots g_{n-1}(g_n(\mathbf{x}; \boldsymbol{\theta}_n), \boldsymbol{\theta}_{n-1}) \dots, \boldsymbol{\theta}_1)$$



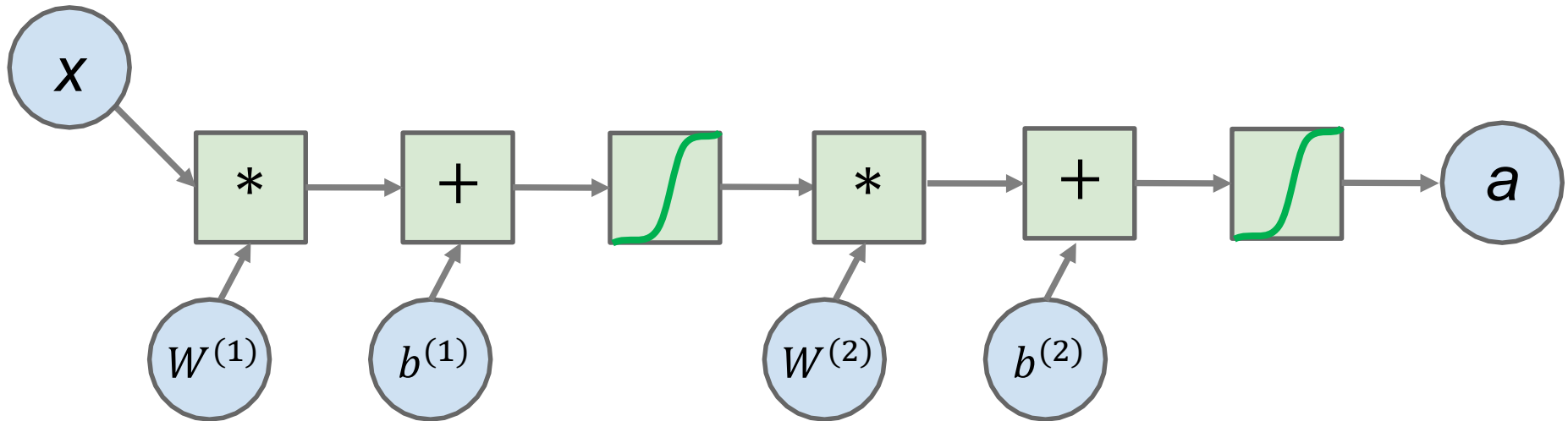
Neural network as variables + operations

- $a = \text{sigmoid}(W^T x + b)$
- Decompose functions into atomic operations
- Separate data (variables) and computing (operations)
- Known as a **computational graph**



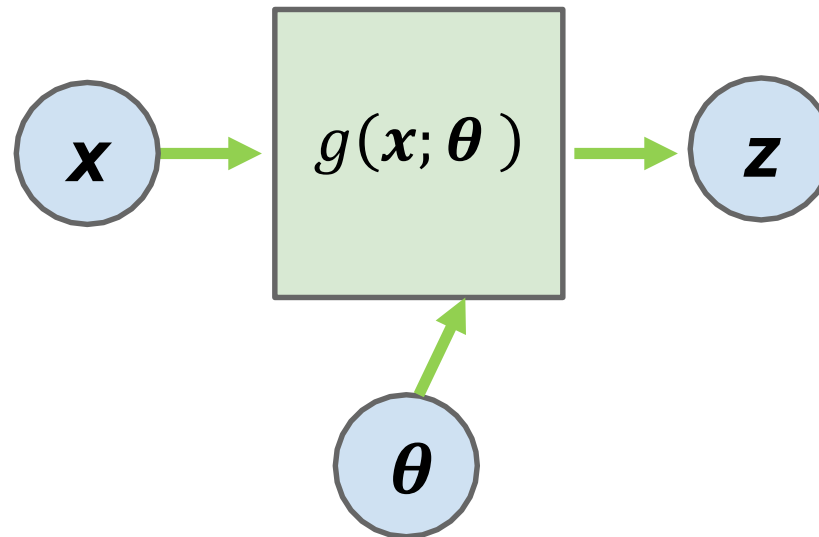
Neural network as variables + operations

- A two-layer neural network



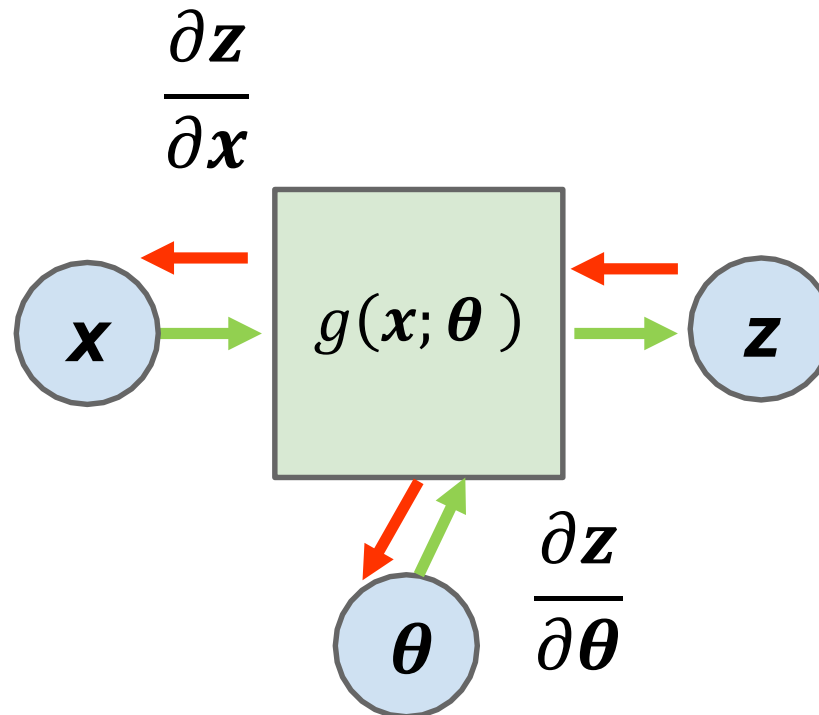
Neural network as variables + operations

- Differentiable operations
- Forward / backward



Neural network as variables + operations

- Differentiable operations
- Forward / backward

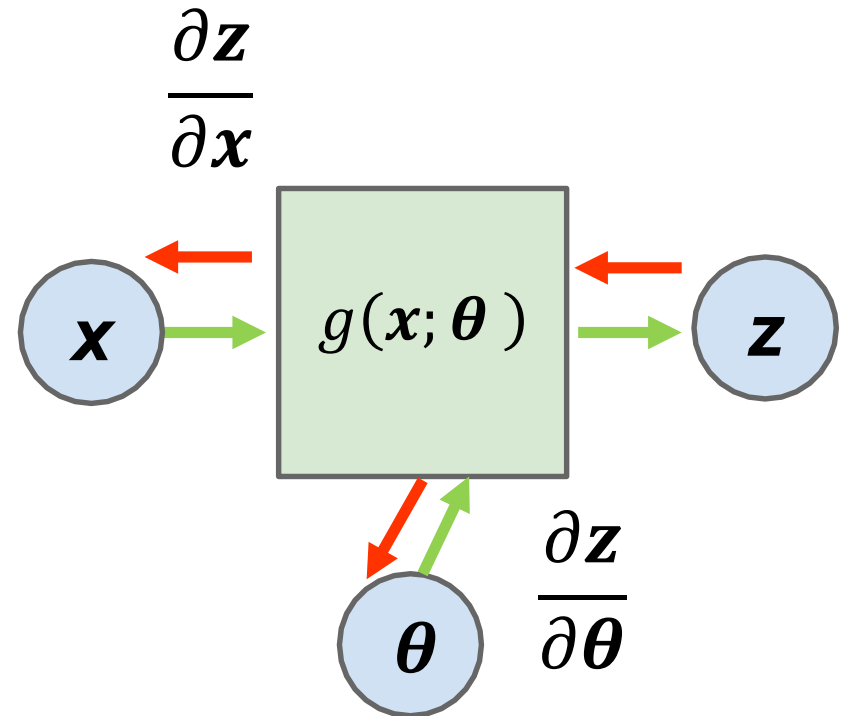


Neural network as variables + operations

- Differentiable operations
- Examples

$$z = g(x) = \textit{sigmoid}(x)$$

$$\frac{\partial z}{\partial x} = z(1 - z)$$



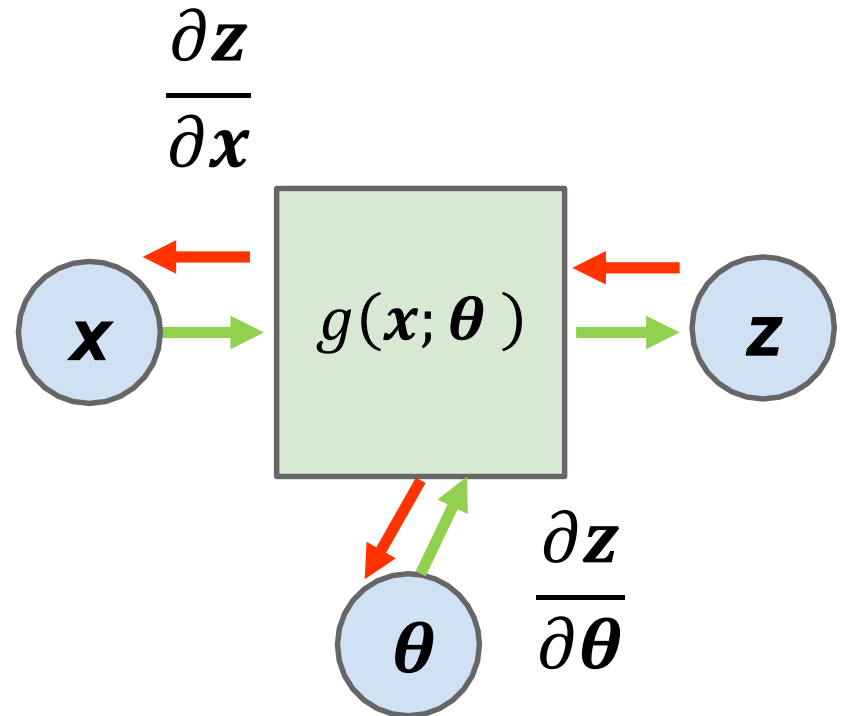
Neural network as variables + operations

- Differentiable operations
- Examples

$$z = g(x; \mathbf{b}) = x + b$$

$$\frac{\partial z}{\partial x} = 1$$

$$\frac{\partial z}{\partial b} = 1$$



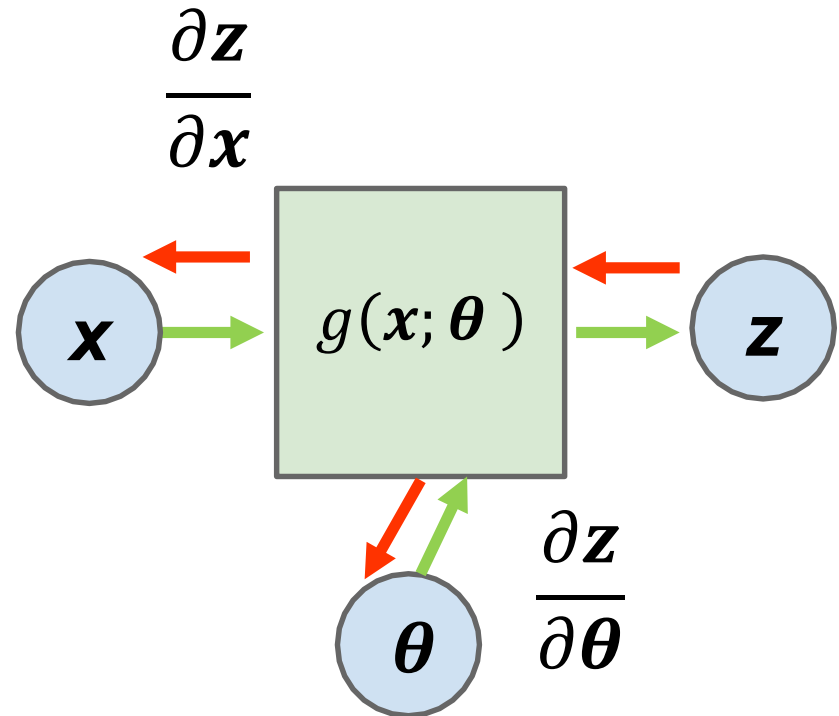
Neural network as variables + operations

- Differentiable operations
- Examples

$$g(\mathbf{x}; \boldsymbol{\theta} = \mathbf{w}) = \mathbf{w}^T \mathbf{x}$$

$$\frac{\partial \mathbf{z}}{\partial \mathbf{x}} = \mathbf{w}$$

$$\frac{\partial \mathbf{z}}{\partial \mathbf{w}} = \mathbf{x}$$



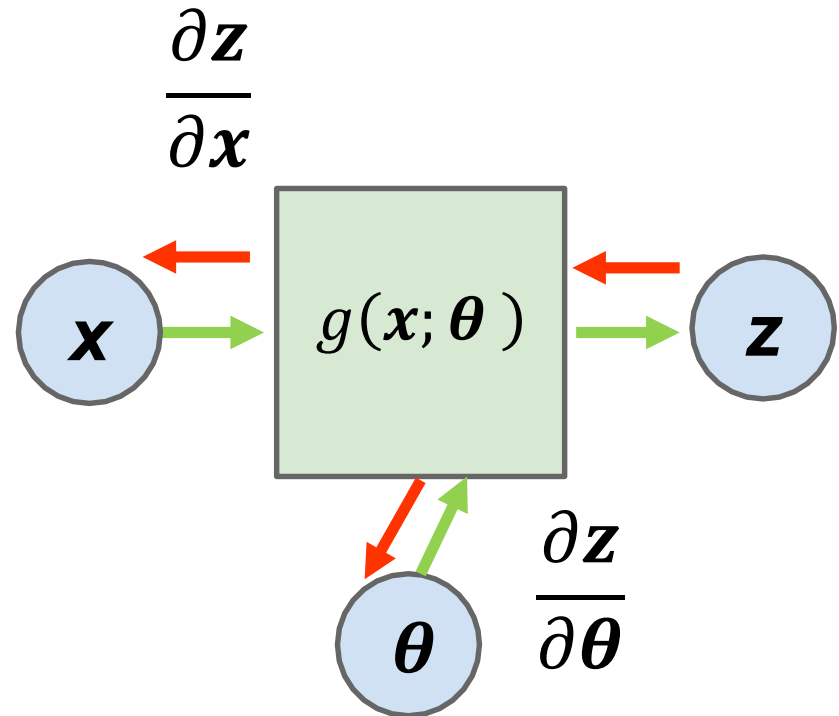
Neural network as variables + operations

- Differentiable operations
- Examples

$$g(x; \theta = W) = W^T x$$

$$\frac{\partial z}{\partial x} = ? \quad \frac{\partial z}{\partial W} = ?$$

$$z = \begin{bmatrix} z_1 \\ z_2 \\ z_3 \end{bmatrix} = \begin{bmatrix} w_1^T \\ w_2^T \\ w_3^T \end{bmatrix} x$$



Neural network as variables + operations

- Differentiable operations
- Examples

$$g(\mathbf{x}; \boldsymbol{\theta} = \mathbf{W}) = \mathbf{W}^T \mathbf{x}$$

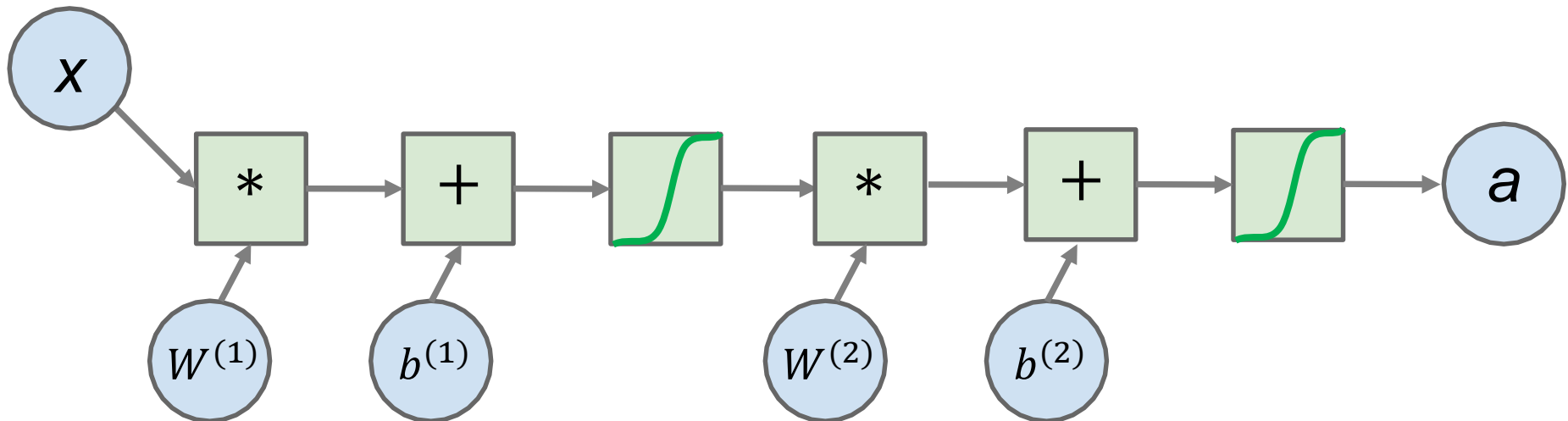
$$\frac{\partial \mathbf{z}}{\partial \mathbf{x}} = ? \quad \frac{\partial \mathbf{z}}{\partial \mathbf{W}} = ?$$

$$\mathbf{z} = \begin{bmatrix} z_1 \\ z_2 \\ z_3 \end{bmatrix} = \begin{bmatrix} \mathbf{w}_1^T \\ \mathbf{w}_2^T \\ \mathbf{w}_3^T \end{bmatrix} \mathbf{x}$$

- Loop over each output value z_i
- Compute $\frac{\partial z_i}{\partial \mathbf{x}} = \mathbf{w}_i$ and $\frac{\partial z_i}{\partial \mathbf{w}_i} = \mathbf{x}$
- Rearrange them into the gradient matrix

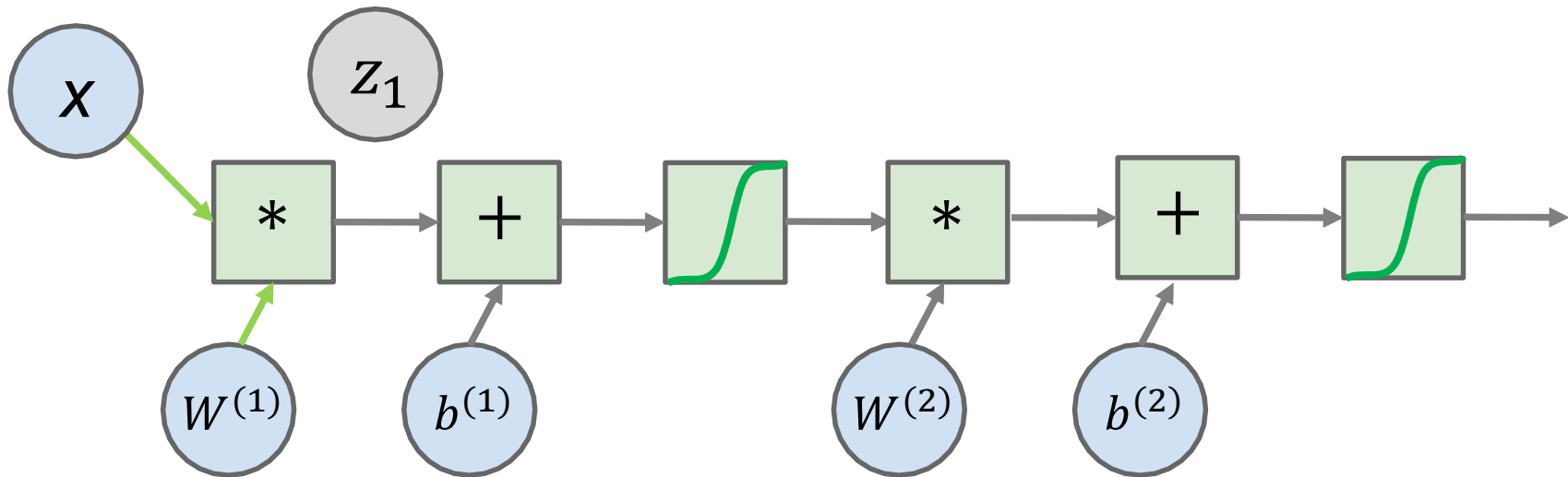
Neural network as a computational graph

- A two-layer neural network
- Forward propagation vs. backward propagation



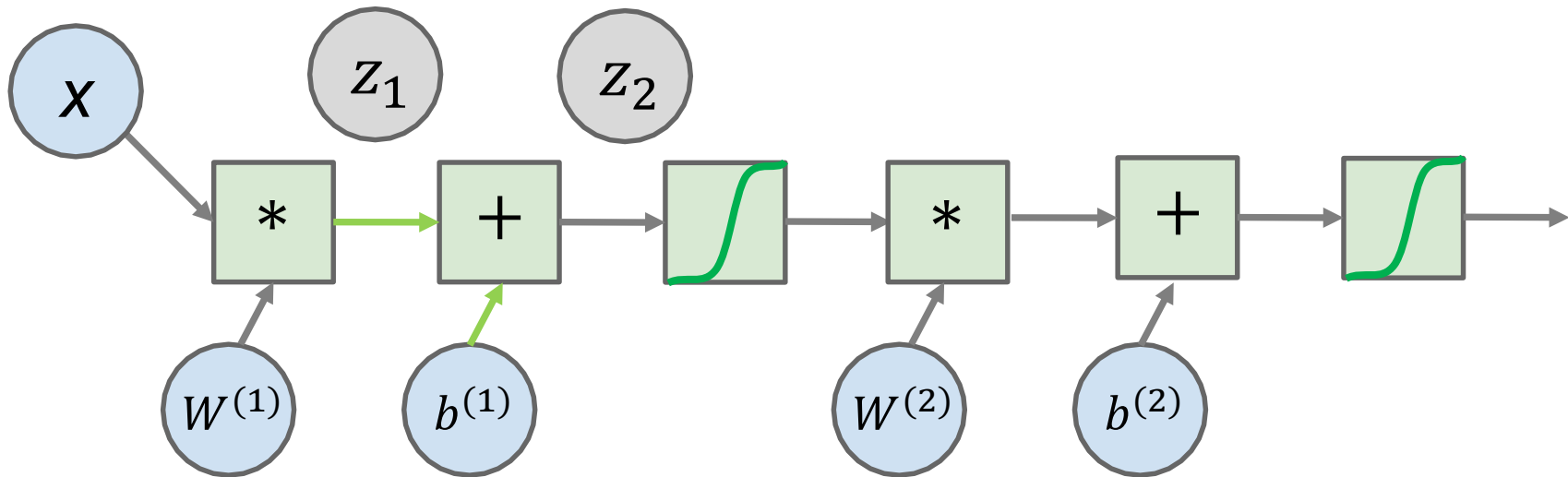
Neural network: forward propagation

- A two-layer neural network
- Intermediate variables Z



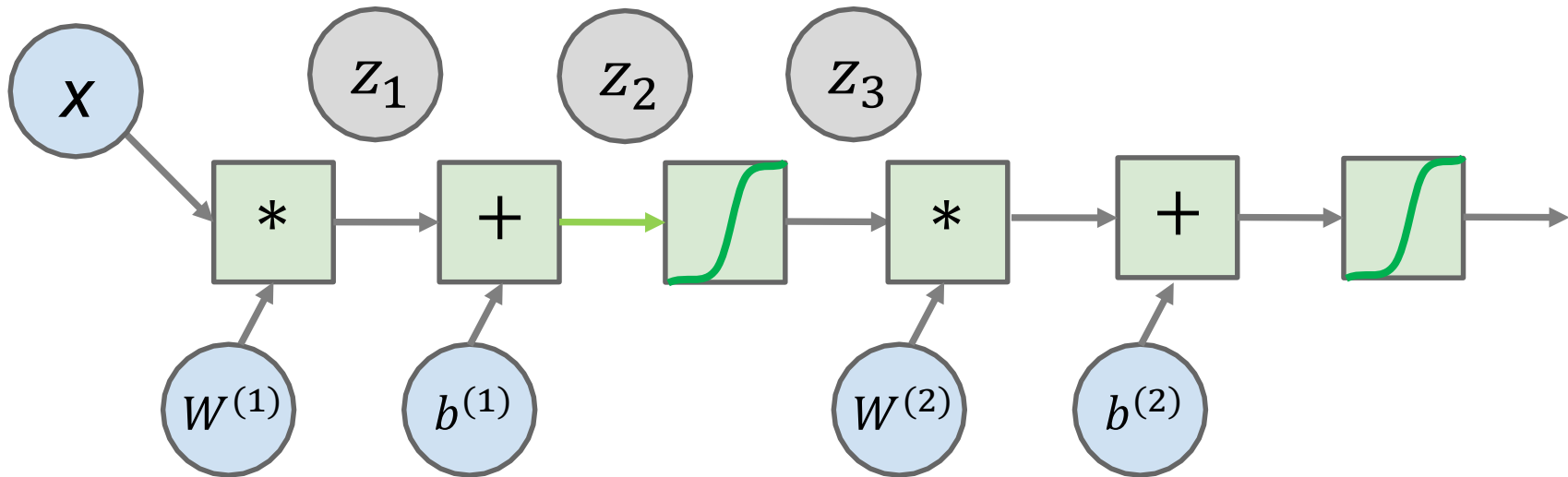
Neural network: forward propagation

- A two-layer neural network
- Intermediate variables Z



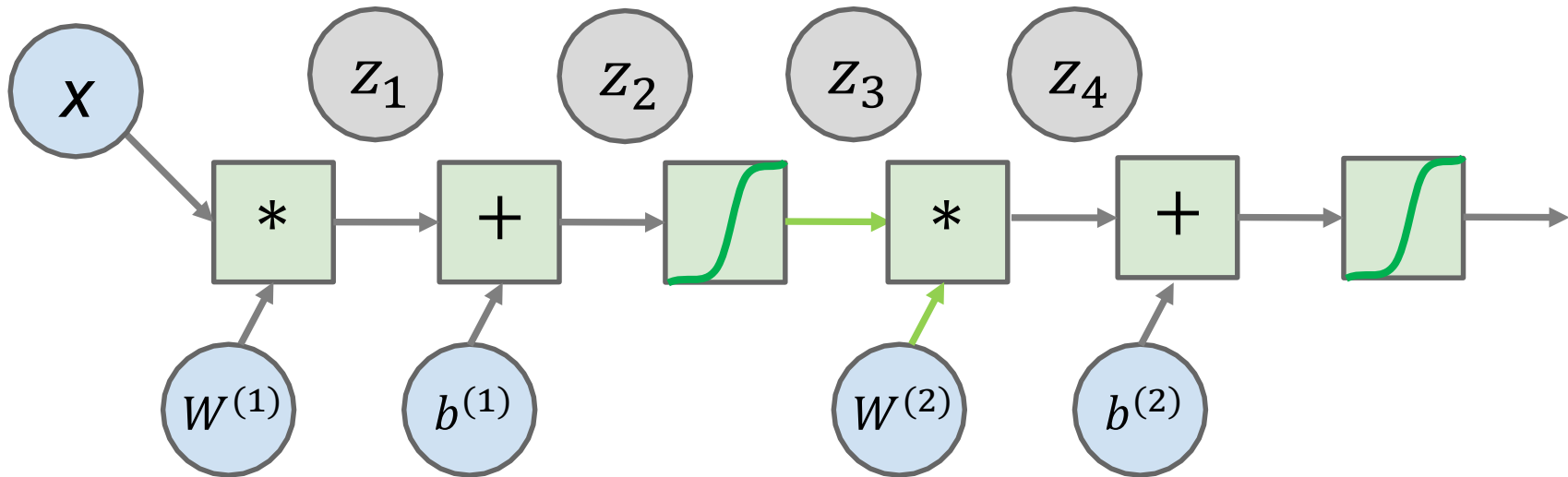
Neural network: forward propagation

- A two-layer neural network
- Intermediate variables Z



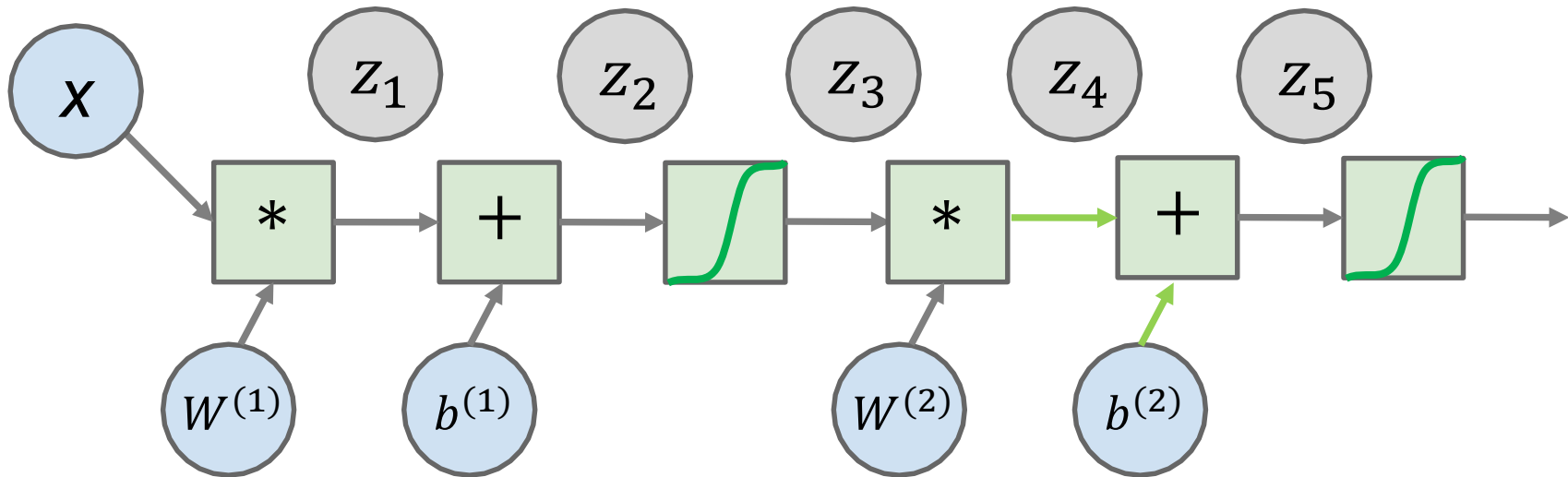
Neural network: forward propagation

- A two-layer neural network
- Intermediate variables Z



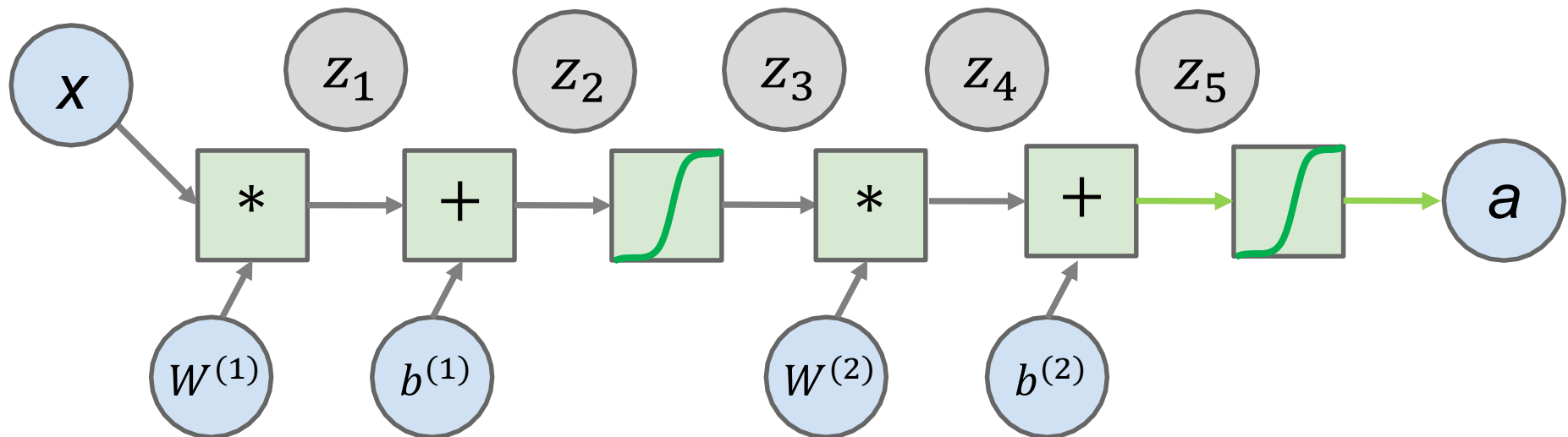
Neural network: forward propagation

- A two-layer neural network
- Intermediate variables Z



Neural network: forward propagation

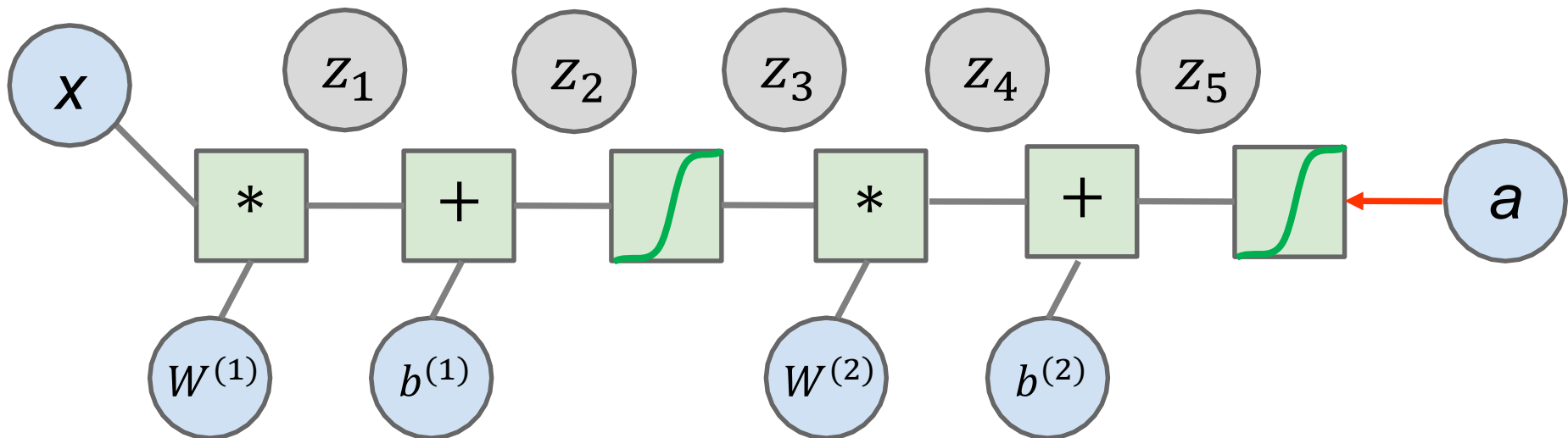
- A two-layer neural network
- Intermediate variables Z
- Forward propagation computes the output a



Neural network: backward propagation

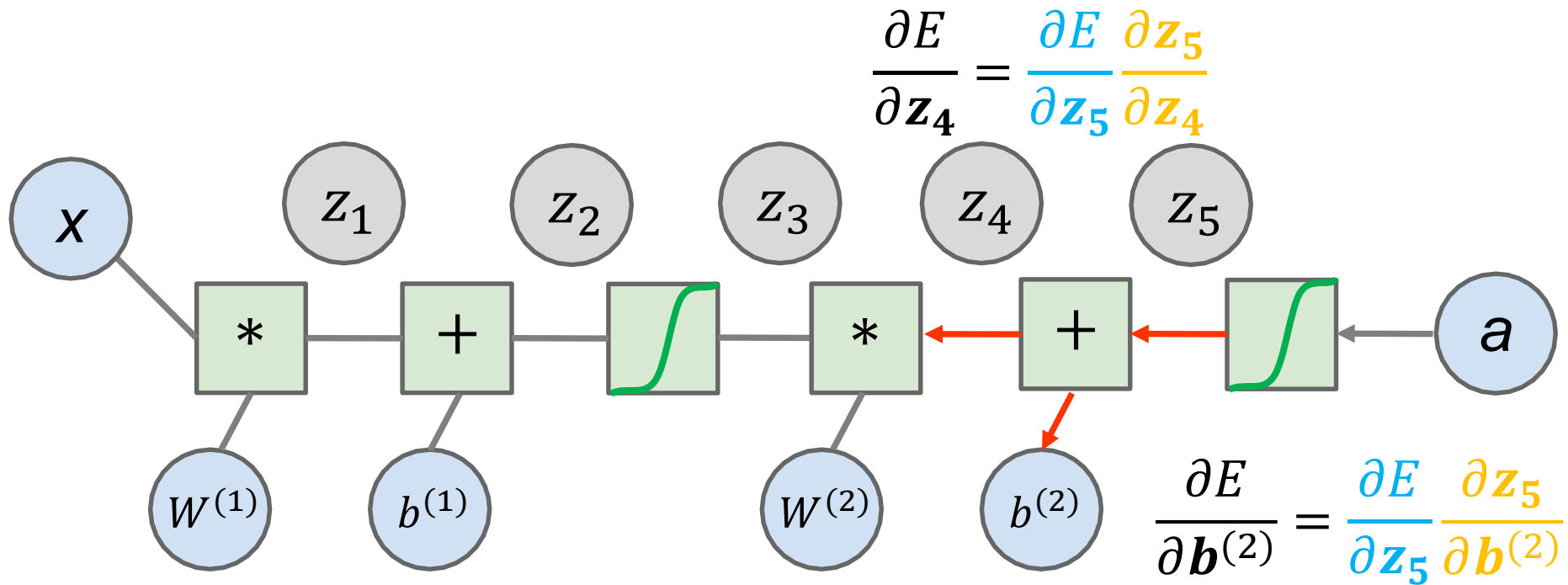
- A two-layer neural network
- Assuming forward propagation is done
- Have all intermediate variables Z
- Minimize a **loss function** E

$$\frac{\partial E}{\partial z_5} = \frac{\partial E}{\partial a} \frac{\partial a}{\partial z_5}$$



Neural network: backward propagation

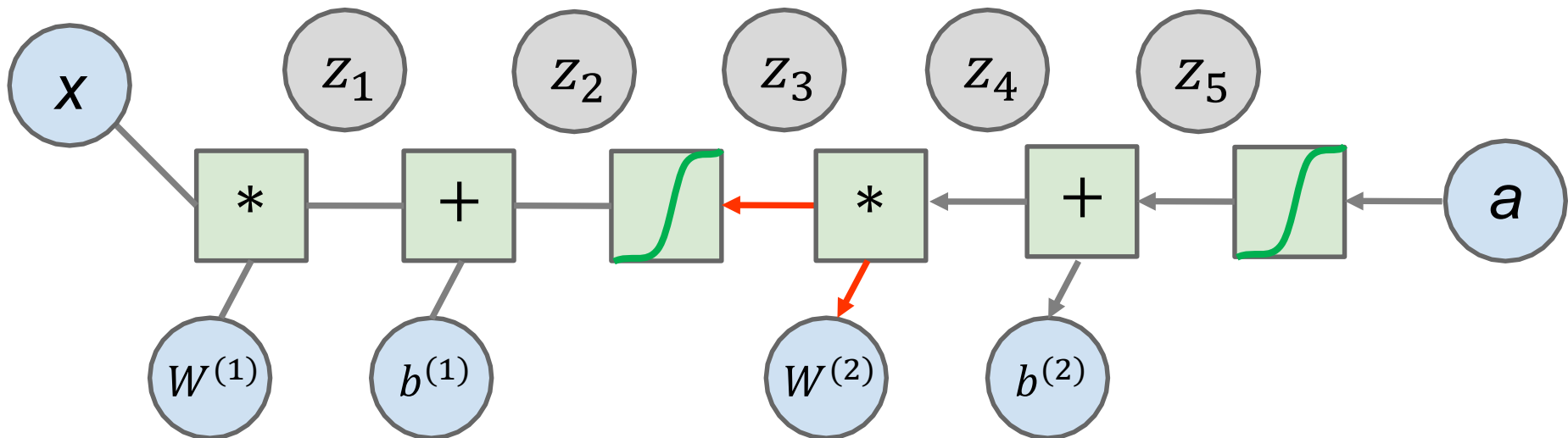
- A two-layer neural network
- Assuming forward propagation is done



Neural network: backward propagation

- A two-layer neural network
- Assuming forward propagation is done

$$\frac{\partial E}{\partial z_3} = \frac{\partial E}{\partial z_4} \frac{\partial z_4}{\partial z_3}$$

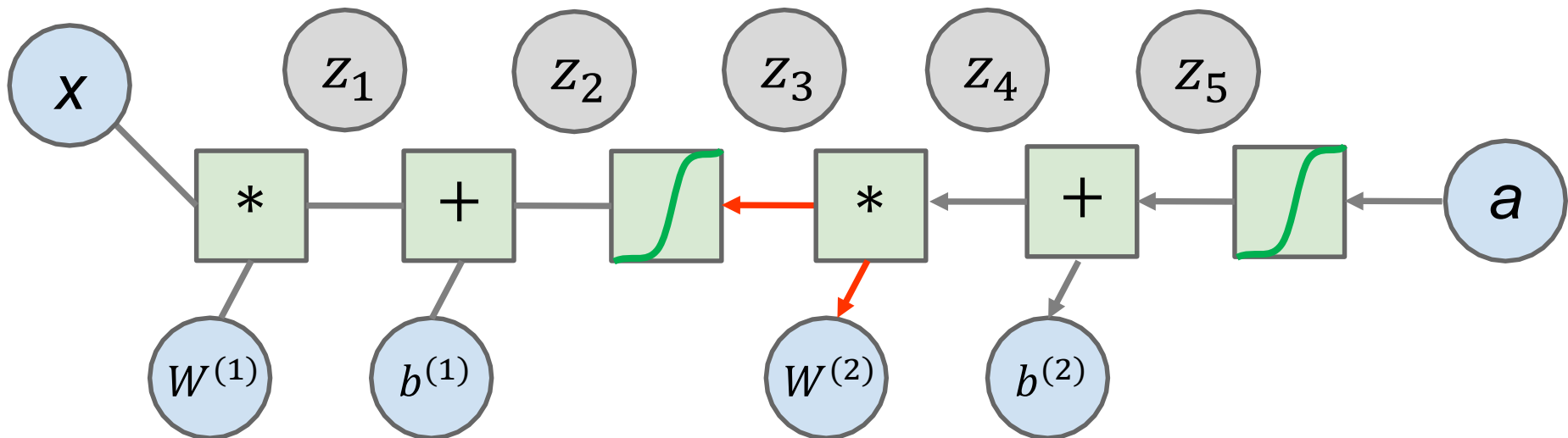


$$\frac{\partial E}{\partial W^{(2)}} = \frac{\partial E}{\partial z_4} \frac{\partial z_4}{\partial W^{(2)}}$$

Neural network: backward propagation

- Gradient to a variable =
gradient on the top x gradient from the current operation

$$\frac{\partial E}{\partial z_3} = \frac{\partial E}{\partial z_4} \frac{\partial z_4}{\partial z_3}$$

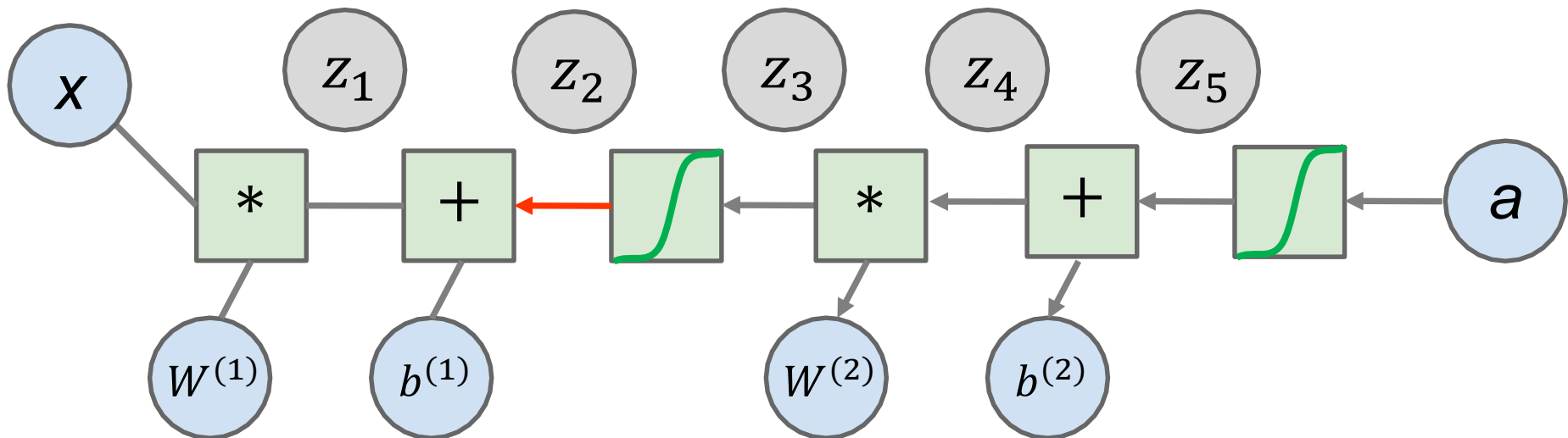


$$\frac{\partial E}{\partial W^{(2)}} = \frac{\partial E}{\partial z_4} \frac{\partial z_4}{\partial W^{(2)}}$$

Neural network: backward propagation

- Gradient to a variable =
gradient on the top x gradient from the current operation

$$\frac{\partial E}{\partial \mathbf{z}_2} = \frac{\partial E}{\partial \mathbf{z}_3} \frac{\partial \mathbf{z}_3}{\partial \mathbf{z}_2}$$



Back propagation: a modern treatment

- Define a neural network as a computational graph
 - Must be a *directed* graph
 - Nodes as *variables* and *operations*
 - All operations must be **differentiable**

Back propagation: a modern treatment

- Define a neural network as a computational graph
- Run topological sort to determine the order of execution of the operations (input -> output)
- **Forward propagation** to compute the output
 - Also store all intermediate outputs
- Compare the output with target label (**the loss function**)
- **Backward propagation** to compute the gradients w.r.t. the model parameters
 - **gradient on the top** x **gradient from the current op**
- Use **gradient descent** to update the parameters