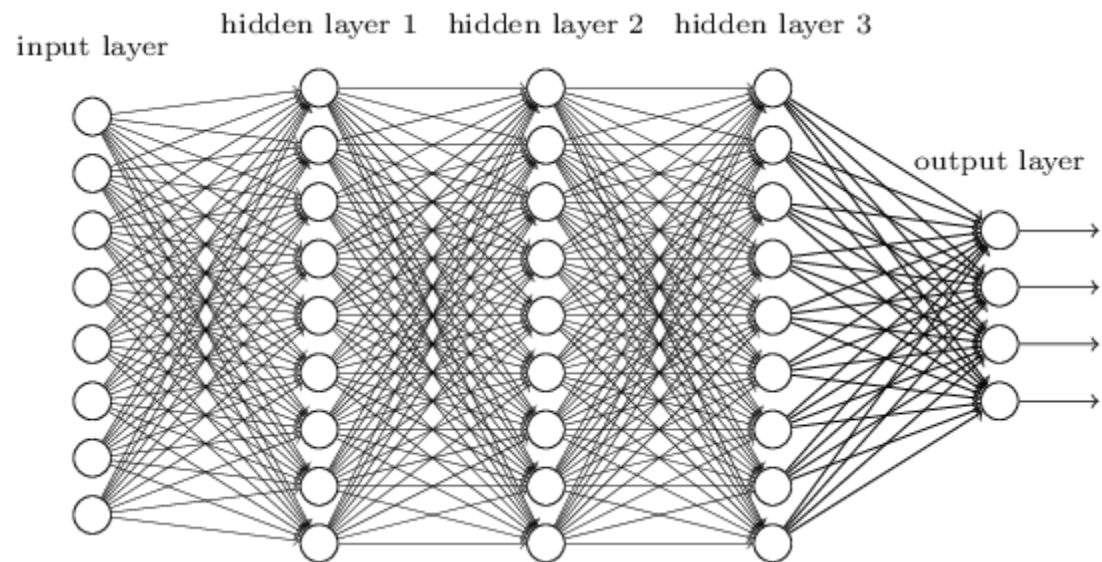


DEEP LEARNING INTRODUCTION



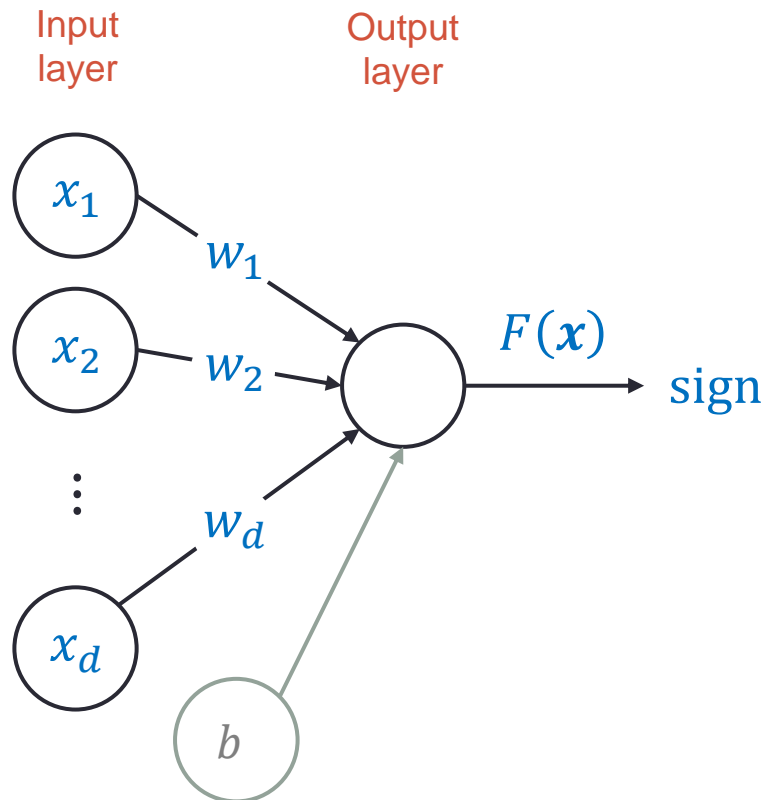
Outline

- From perceptrons to neural networks
- Training neural networks
- Backpropagation
- Playground

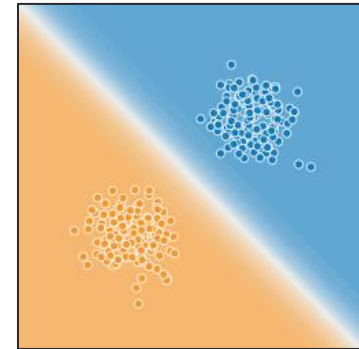
Perceptron

- Recall the perceptron linear model

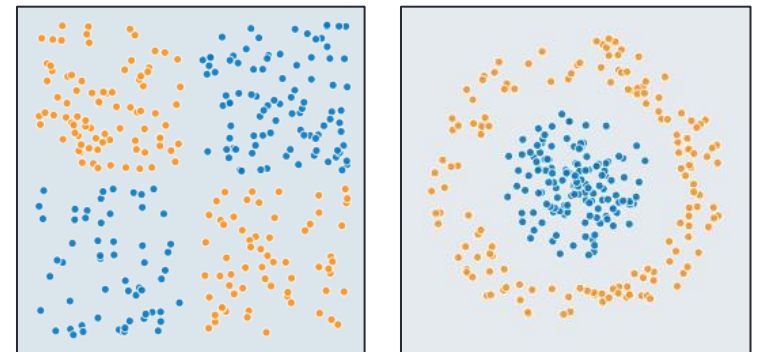
$$F(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b = \sum_{i=1}^d x_i w_i + b$$



- Creates a linear decision boundary.

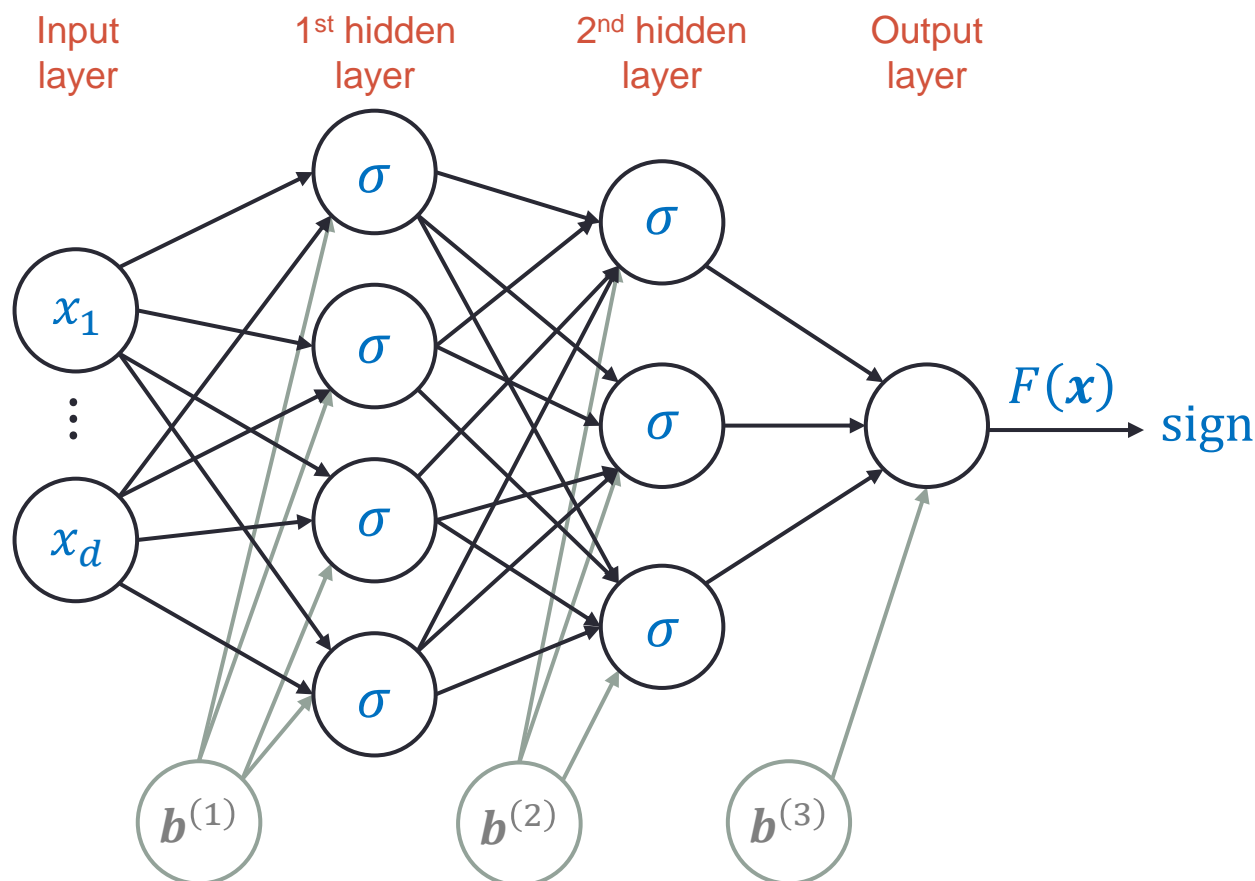


- How can we create more complicated decision boundaries?

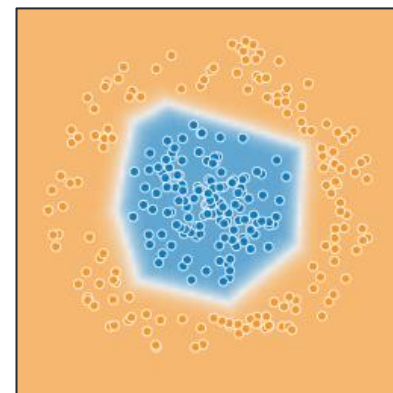
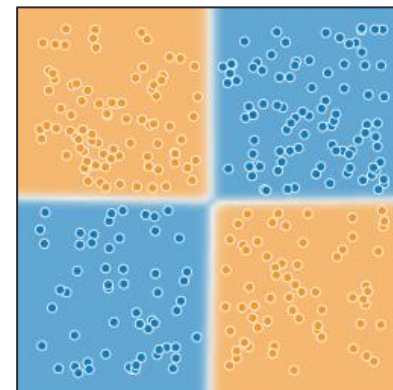


Neural networks (Multilayer Perceptron)

- We can add more layers!
- For instance, $F(\mathbf{x}) = \mathbf{w}^{(3)\top} \sigma \left(\mathbf{W}^{(2)\top} \sigma \left(\mathbf{W}^{(1)\top} \mathbf{x} + \mathbf{b}^{(1)} \right) + \mathbf{b}^{(2)} \right) + \mathbf{b}^{(3)}$

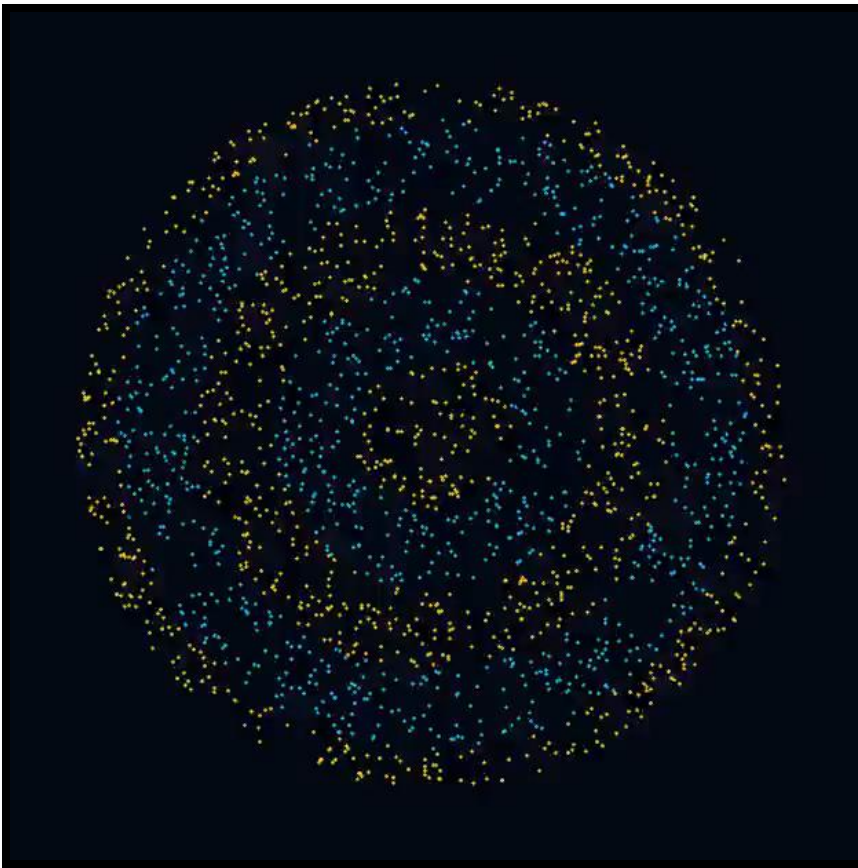


- Can now fit richer functions

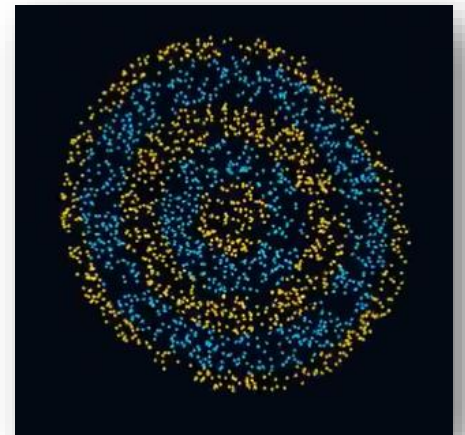


Demo: “Feature” learning

- Two-layer neural network tries to separate “circles” by learning a single new feature.



Source: [Matt Henderson on Twitter](#)



Training neural networks

- Define a network F parameterized by $\Theta = (\mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \mathbf{W}^{(2)}, \mathbf{b}^{(2)}, \dots)$. *Which?*
- Define a loss over network outputs, e.g., MSE: $\mathcal{L} = \frac{1}{m} \sum_{i=1}^m (F(\mathbf{x}_i) - y_i)^2$ *Which?*
- **Gradient descent** general scheme:
 1. Initialize parameters randomly
 2. While model has not converged:
 - i. Compute gradient $\frac{\partial}{\partial \Theta} \mathcal{L}(\Theta)$ *How?*
 - ii. Update weights $\Theta \leftarrow \Theta - \eta \frac{\partial}{\partial \Theta} \mathcal{L}(\Theta)$ *How?*
- The loss is usually not convex w.r.t Θ ,
but practically we often converge to a good minimum!

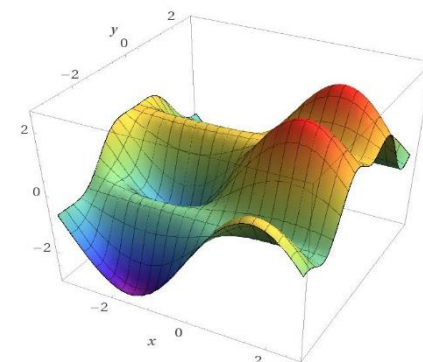
How?

When?

How?

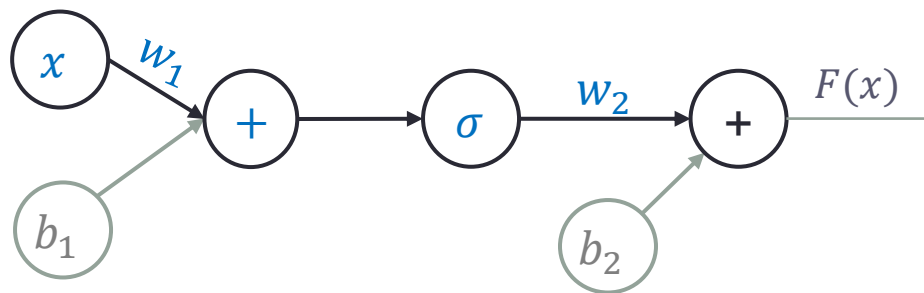
How?

Today's focus



Simple regression example

- Consider a simple network $F(x) = w_2 \cdot \sigma(w_1x + b_1) + b_2$.



- The architecture is defined only after choosing the activation function.

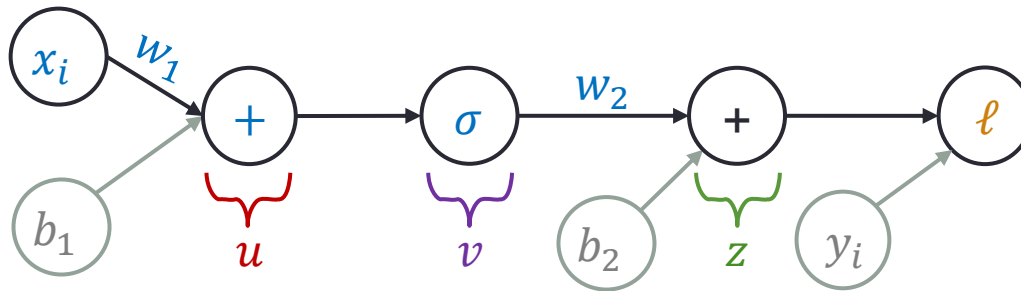
We will use the sigmoid function $\sigma(a) = \frac{1}{1+e^{-a}}$.

- Trying to solve a regression problem, we minimize the squared loss

$$\mathcal{L} = \sum_{i=1}^m (F(x_i) - y_i)^2 = \sum_{i=1}^m \ell(x_i, y_i)$$

Simple regression example

- Consider a simple network $F(x) = w_2 \cdot \sigma(w_1x + b_1) + b_2$.



where $\sigma(a) = \frac{1}{1+e^{-a}}$ and the loss is $\mathcal{L} = \sum_{i=1}^m (F(x_i) - y_i)^2 = \sum_{i=1}^m \ell(x_i, y_i)$.

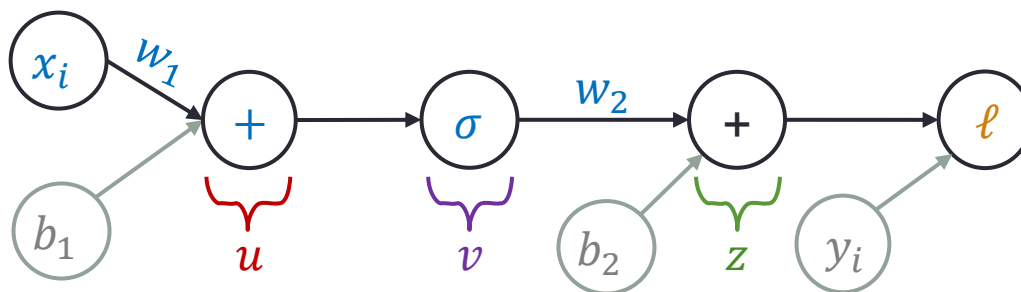
- Denote: $u = w_1x + b_1$, $v = \sigma(u)$, $z = w_2v + b_2 (= F(x))$, $\ell = (z - y_i)^2$
- Compute all partial derivatives using the **chain rule**:

$$\frac{\partial}{\partial w_2} \ell = \frac{\partial \ell}{\partial z} \frac{\partial z}{\partial w_2} =$$

$$\frac{\partial}{\partial b_2} \ell = \frac{\partial \ell}{\partial z} \frac{\partial z}{\partial b_2} =$$

Simple regression example

- Consider a simple network $F(x) = w_2 \cdot \sigma(w_1x + b_1) + b_2$.



where $\sigma(a) = \frac{1}{1+e^{-a}}$ and the loss is $\mathcal{L} = \sum_{i=1}^m (F(x_i) - y_i)^2 = \sum_{i=1}^m \ell(x_i, y_i)$.

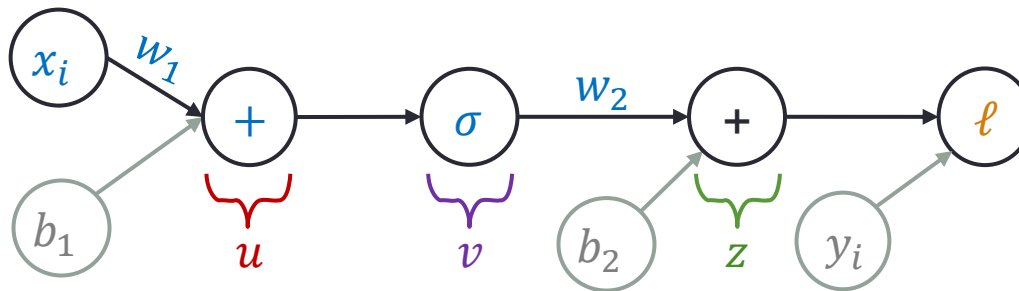
- Denote: $u = w_1x + b_1$, $v = \sigma(u)$, $z = w_2v + b_2 (= F(x))$, $\ell = (z - y_i)^2$
- Compute all partial derivatives using the **chain rule**:

$$\frac{\partial}{\partial w_2} \ell = \frac{\partial \ell}{\partial z} \frac{\partial z}{\partial w_2} = 2(z - y_i) \cdot v$$

$$\frac{\partial}{\partial b_2} \ell = \frac{\partial \ell}{\partial z} \frac{\partial z}{\partial b_2} = 2(z - y_i) \cdot 1$$

Simple regression example

- Consider a simple network $F(x) = w_2 \cdot \sigma(w_1x + b_1) + b_2$.



where $\sigma(a) = \frac{1}{1+e^{-a}}$ and the loss is $\mathcal{L} = \sum_{i=1}^m (F(x_i) - y_i)^2 = \sum_{i=1}^m \ell(x_i, y_i)$.

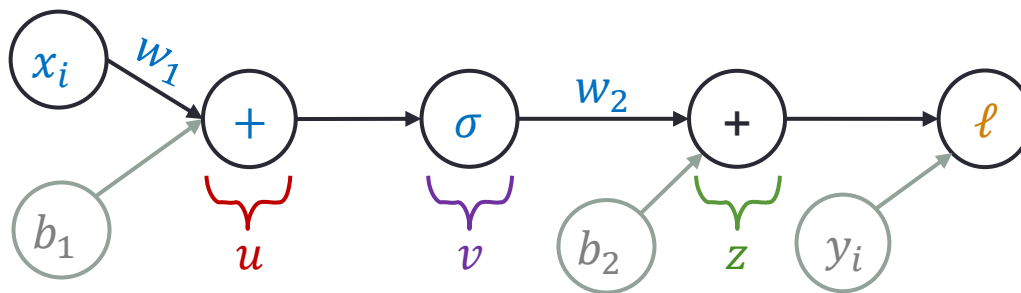
- Denote: $u = w_1x + b_1$, $v = \sigma(u)$, $z = w_2v + b_2 (= F(x))$, $\ell = (z - y_i)^2$
- Compute all partial derivatives using the **chain rule**:

$$\frac{\partial}{\partial w_1} \ell = \frac{\partial \ell}{\partial z} \frac{\partial z}{\partial v} \frac{\partial v}{\partial u} \frac{\partial u}{\partial w_1} =$$

$$\frac{\partial}{\partial b_1} \ell = \frac{\partial \ell}{\partial z} \frac{\partial z}{\partial v} \frac{\partial v}{\partial u} \frac{\partial u}{\partial b_1} =$$

Simple regression example

- Consider a simple network $F(x) = w_2 \cdot \sigma(w_1 x + b_1) + b_2$.



where $\sigma(a) = \frac{1}{1+e^{-a}}$ and the loss is $\mathcal{L} = \sum_{i=1}^m (F(x_i) - y_i)^2 = \sum_{i=1}^m \ell(x_i, y_i)$.

- Denote: $u = w_1 x + b_1$, $v = \sigma(u)$, $z = w_2 v + b_2 (= F(x))$, $\ell = (z - y_i)^2$

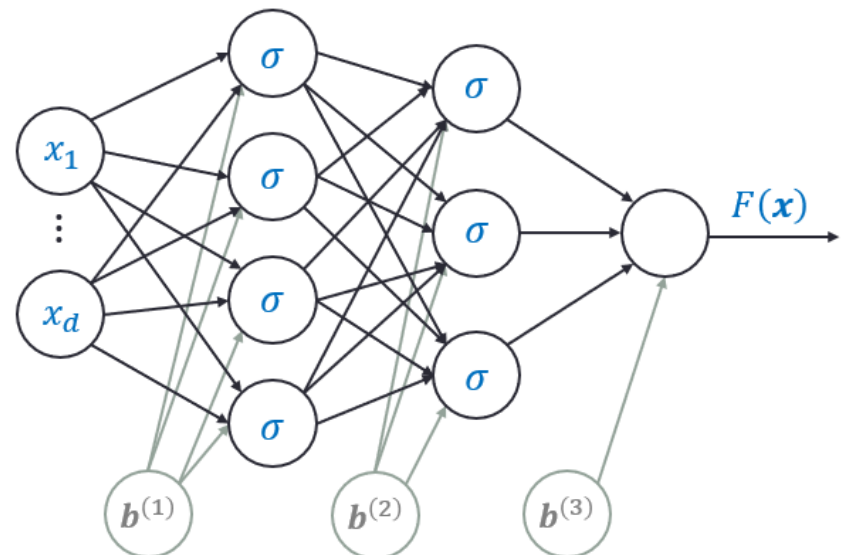
- Compute all partial derivatives using the **chain rule**: $\frac{\ell(x)}{f(x)} = \frac{f(x) \cdot g(x) - g'(x) \cdot f(x)}{g^2(x)}$

$$\frac{\partial}{\partial w_1} \ell = \frac{\partial \ell}{\partial z} \frac{\partial z}{\partial v} \frac{\partial v}{\partial u} \frac{\partial u}{\partial w_1} = 2(z - y_i) \cdot w_2 \cdot \sigma(u) (1 - \sigma(u)) \cdot x$$

$$\frac{\partial}{\partial b_1} \ell = \frac{\partial \ell}{\partial z} \frac{\partial z}{\partial v} \frac{\partial v}{\partial u} \frac{\partial u}{\partial b_1} = 2(z - y_i) \cdot w_2 \cdot \sigma(u) (1 - \sigma(u)) \cdot 1$$

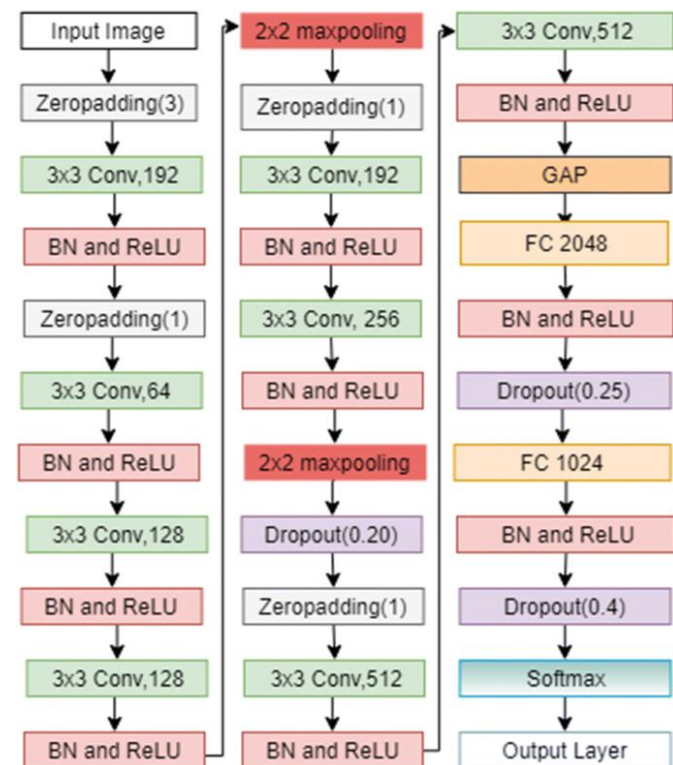
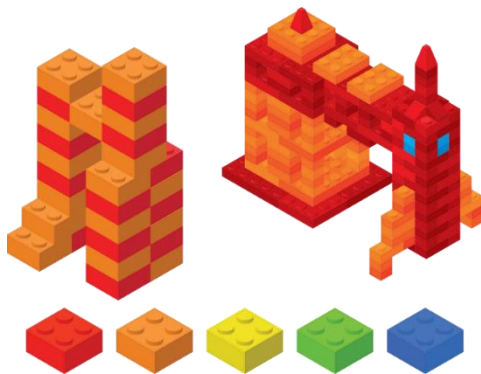
Larger models

- Same idea!
- Compute **all partial derivatives** of the loss w.r.t the parameters.
 - For a large network, this is time consuming.
 - Backpropagation does this efficiently!
 - We will present the modular approach.



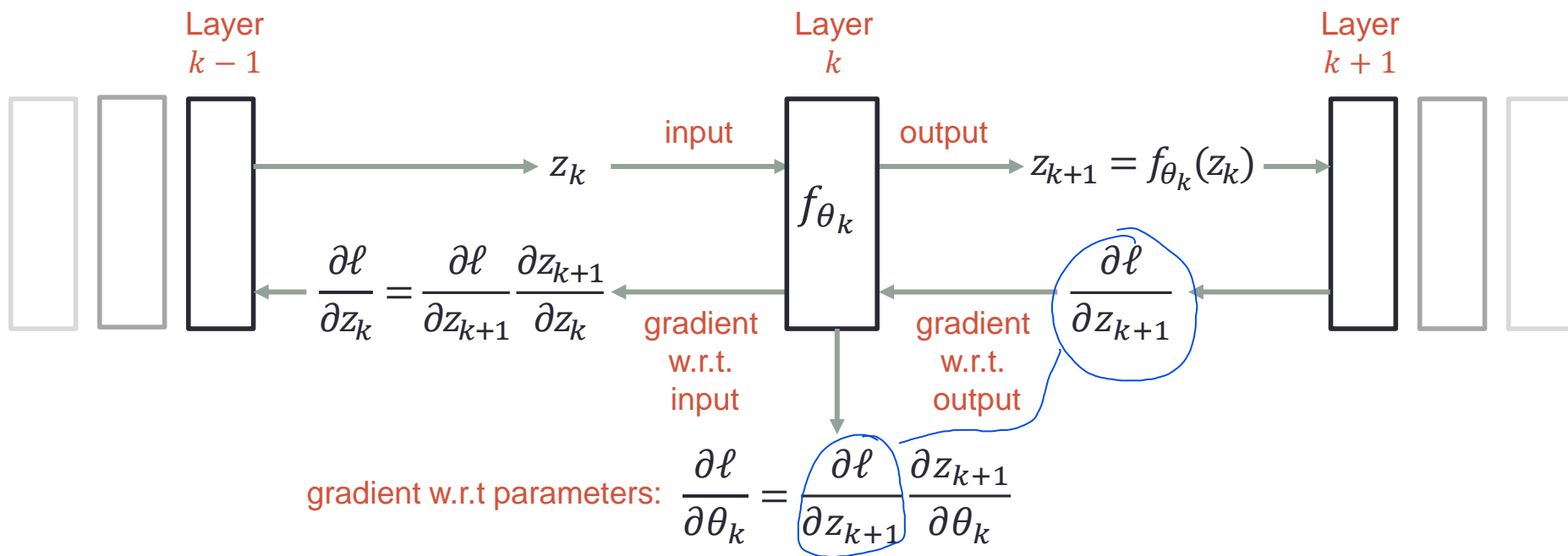
Modular Approach: Code layers, not networks!

- **Layer specification** – each layer should provide three functions:
 1. The layer output given its input (Forward)
 2. Loss gradient w.r.t. the input (Backward)
 3. Loss gradient w.r.t. parameters

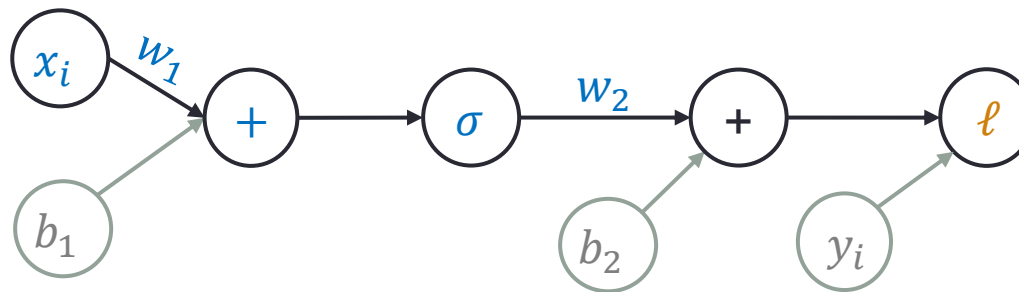


Modular Approach: Code layers not networks!

- Layer specification – each layer should provide three functions:
 1. The layer output given its input (Forward)
 2. Loss gradient w.r.t. the input (Backward)
 3. Loss gradient w.r.t. parameters
- For instance, layer k implements a function f_{θ_k} with its parameters θ_k .



- Let's build the regression network we saw earlier with the modular approach



Linear

$$f(z) = w_1 z + b_1$$

$$\frac{\partial f}{\partial z} = w_1$$

$$\frac{\partial f}{\partial w_1} = z$$

$$\frac{\partial f}{\partial b_1} = 1$$

ReLU

$$f(z) = \max\{0, z\}$$

$$\frac{\partial f}{\partial z} = \mathbb{I}_{z>0}$$

Linear

$$f(z) = w_2 z + b_2$$

$$\frac{\partial f}{\partial z} = w_2$$

$$\frac{\partial f}{\partial w_2} = z$$

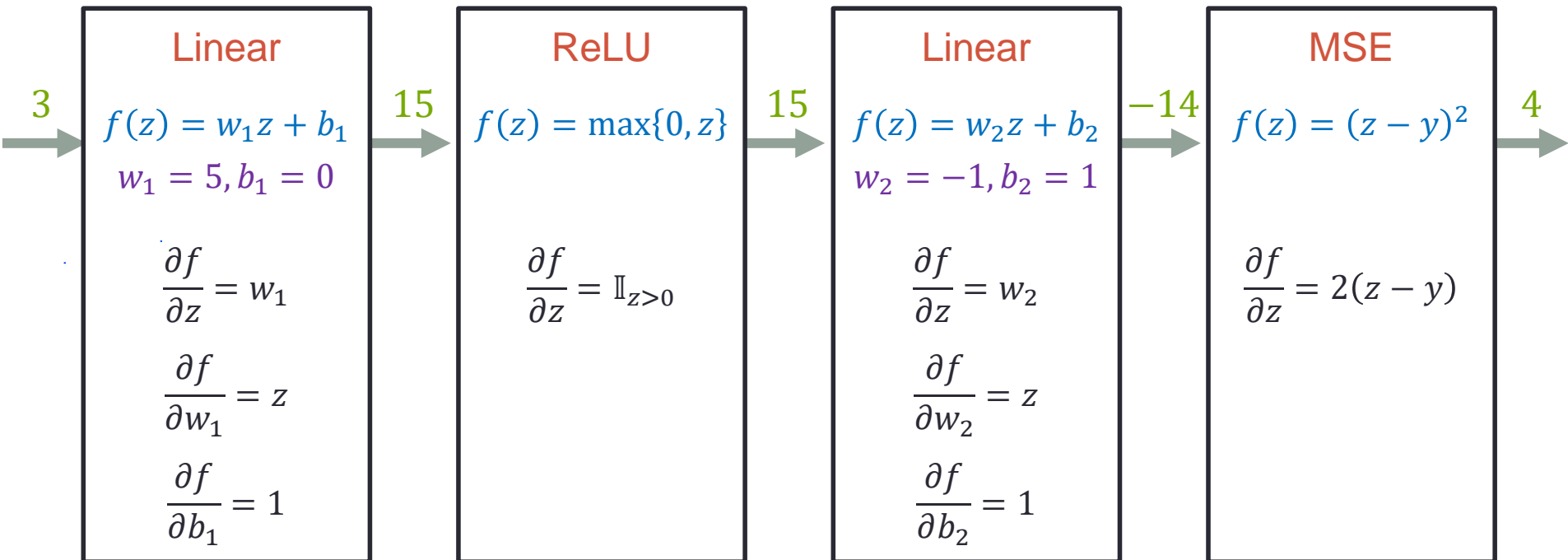
$$\frac{\partial f}{\partial b_2} = 1$$

MSE

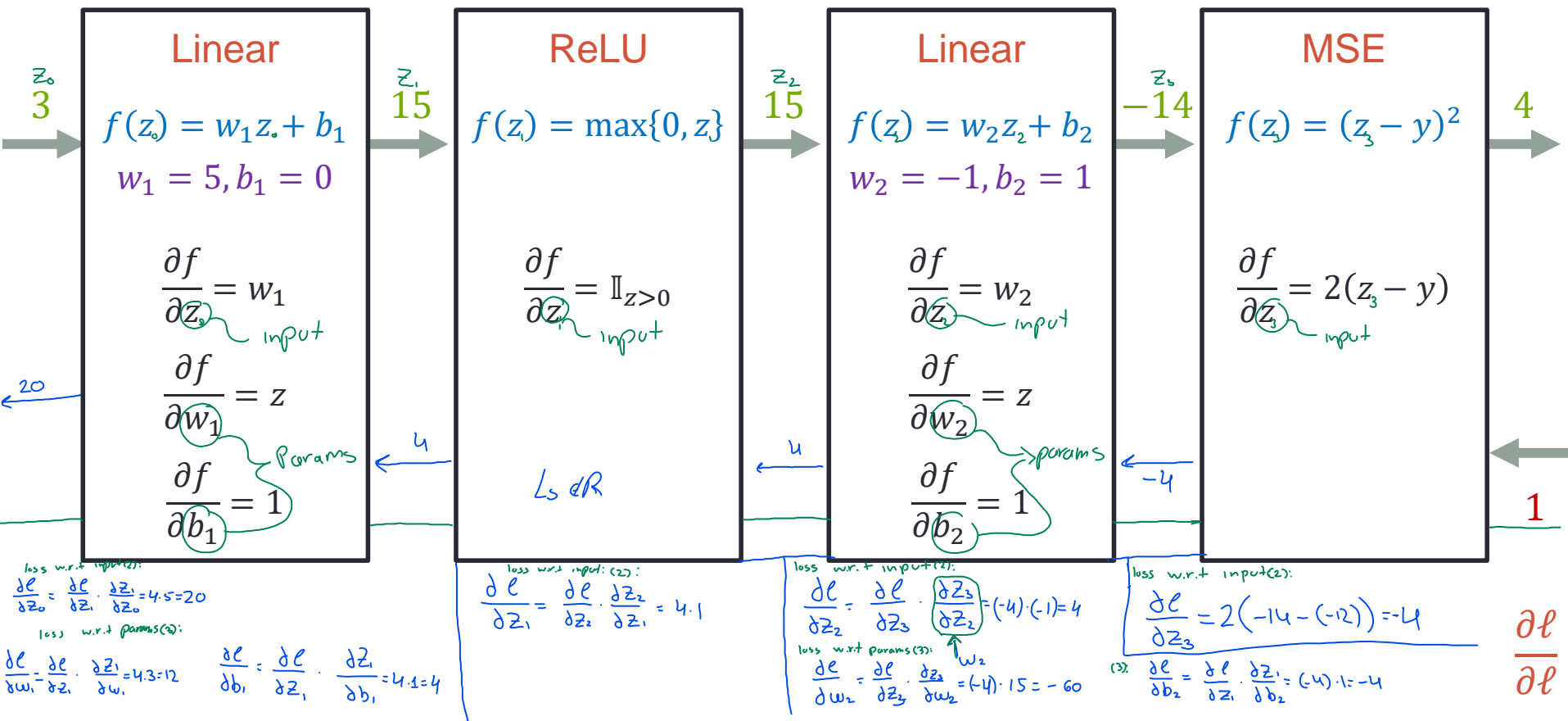
$$f(z) = (z - y)^2$$

$$\frac{\partial f}{\partial z} = 2(z - y)$$

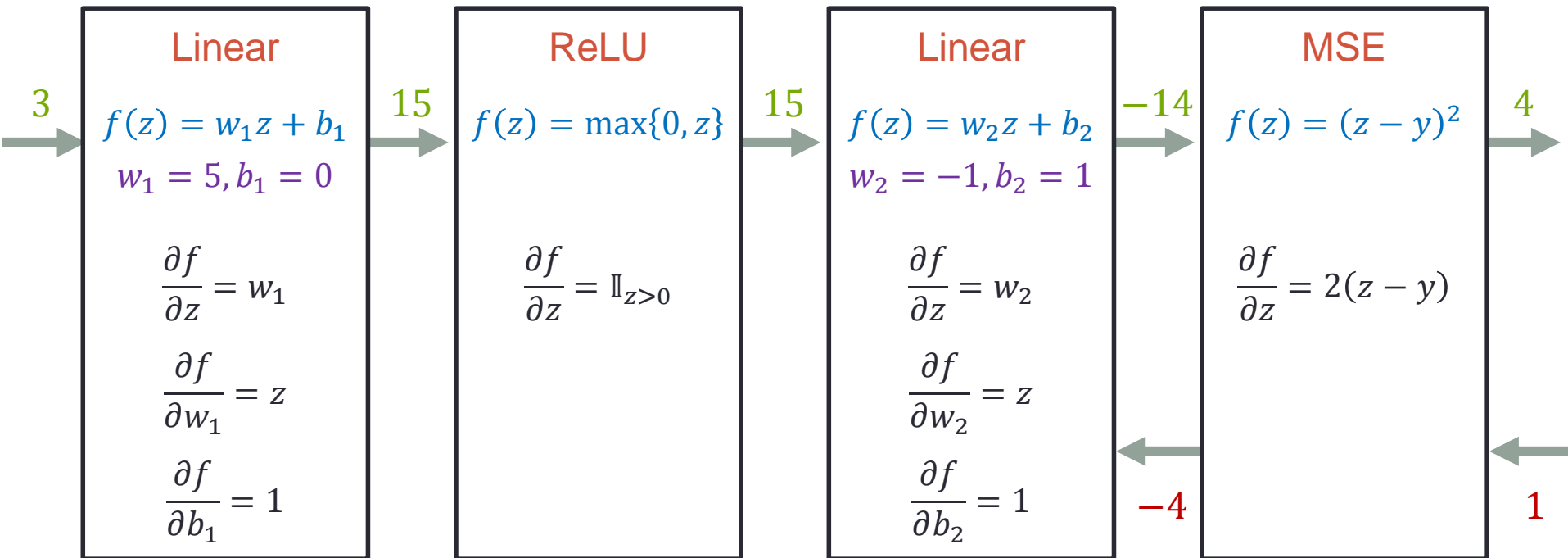
- Initialization: $w_1 = 5, w_2 = -1, b_1 = 0, b_2 = 1$
- Input: $x = 3, y = -12$
- First, we compute the **forward** pass.



- Initialization: $w_1 = 5, w_2 = -1, b_1 = 0, b_2 = 1$
- Input: $x = 3, y = -12$
- First, we compute the **forward** pass.
- Then, to compute gradients w.r.t θ , we perform a **backward** pass.

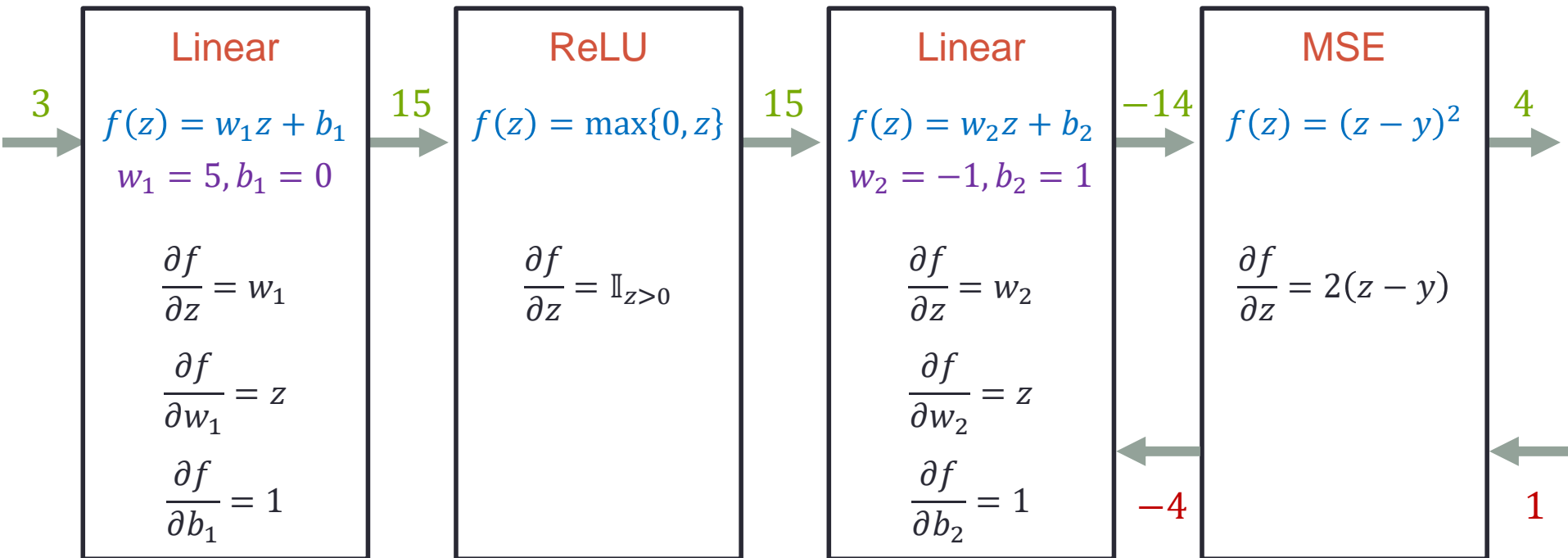


- Initialization: $w_1 = 5, w_2 = -1, b_1 = 0, b_2 = 1$
- Input: $x = 3, y = -12$
- First, we compute the **forward** pass.
- Then, to compute gradients w.r.t θ , we perform a **backward** pass.



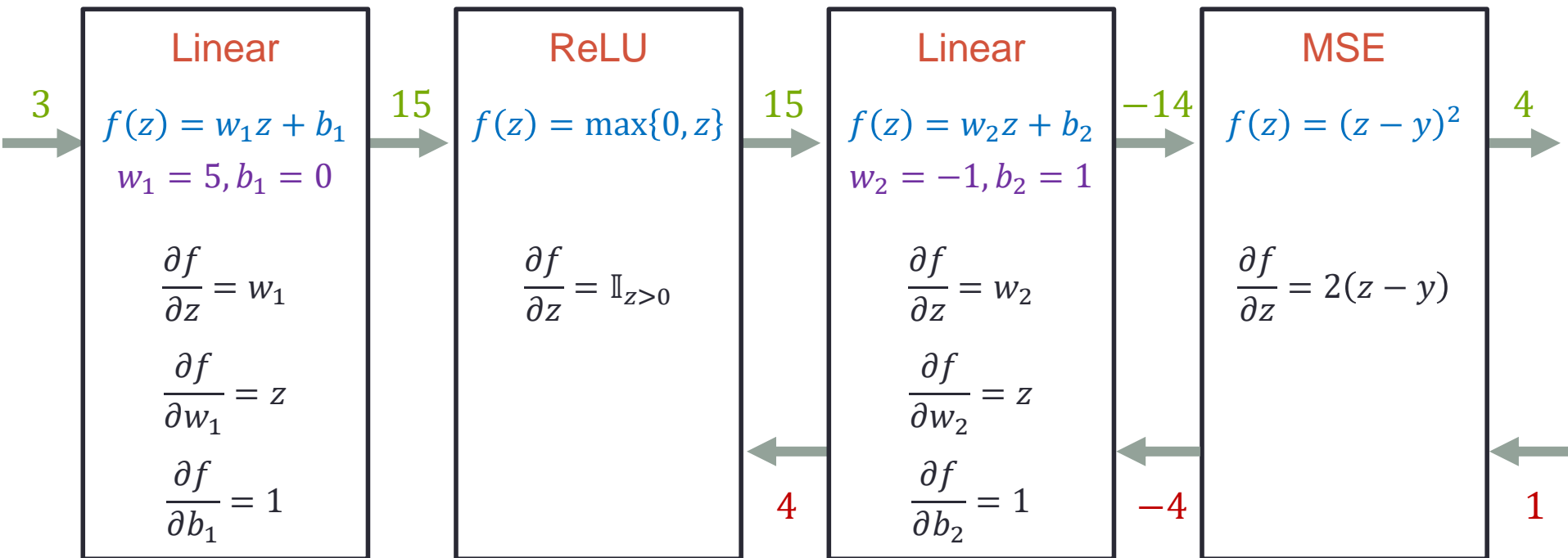
$$\frac{\partial \ell}{\partial z} = \frac{\partial \ell}{\partial f} \frac{\partial f}{\partial z} = 1 \cdot -4$$

- Initialization: $w_1 = 5, w_2 = -1, b_1 = 0, b_2 = 1$
- Input: $x = 3, y = -12$
- First, we compute the **forward** pass.
- Then, to compute gradients w.r.t θ , we perform a **backward** pass.



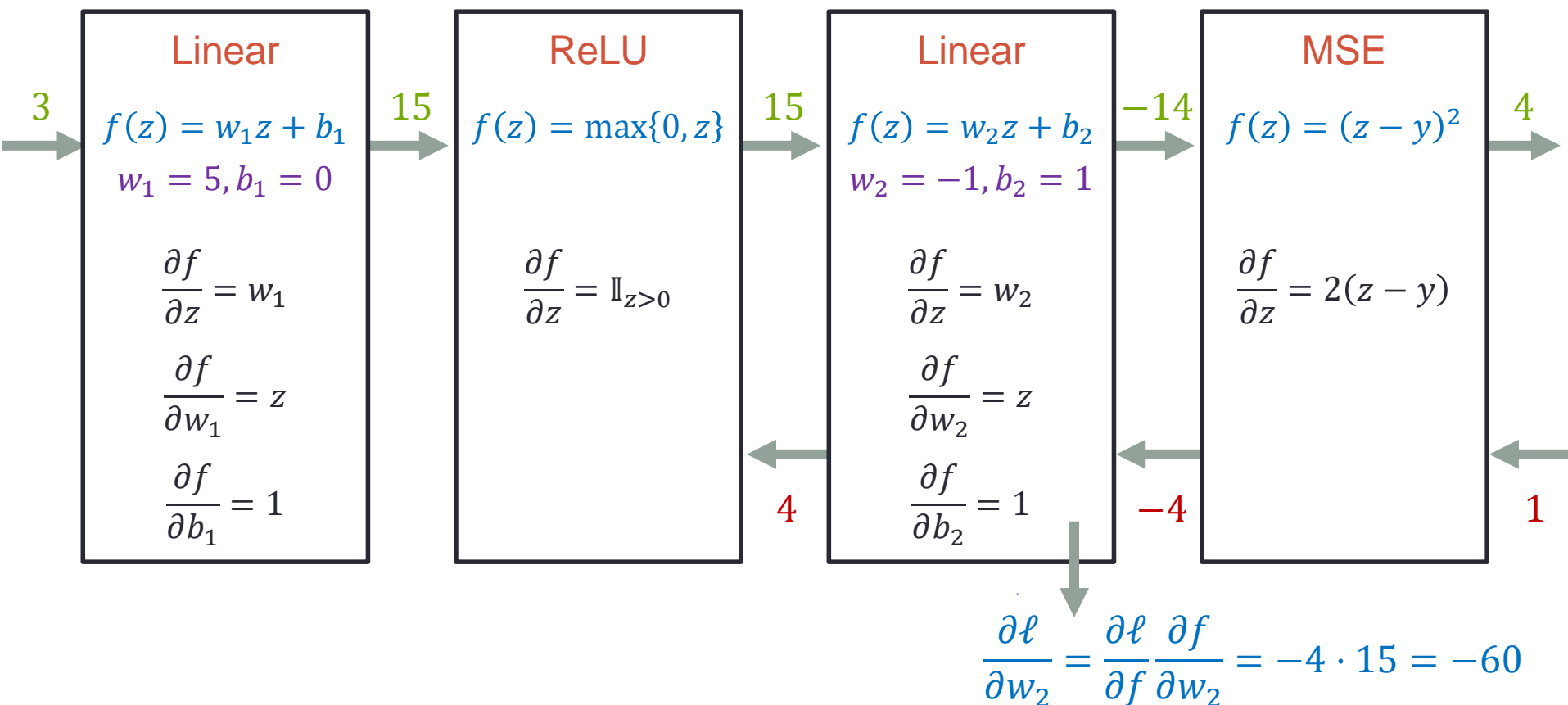
(no parameters for this layer)

- Initialization: $w_1 = 5, w_2 = -1, b_1 = 0, b_2 = 1$
- Input: $x = 3, y = -12$
- First, we compute the **forward** pass.
- Then, to compute gradients w.r.t θ , we perform a **backward** pass.

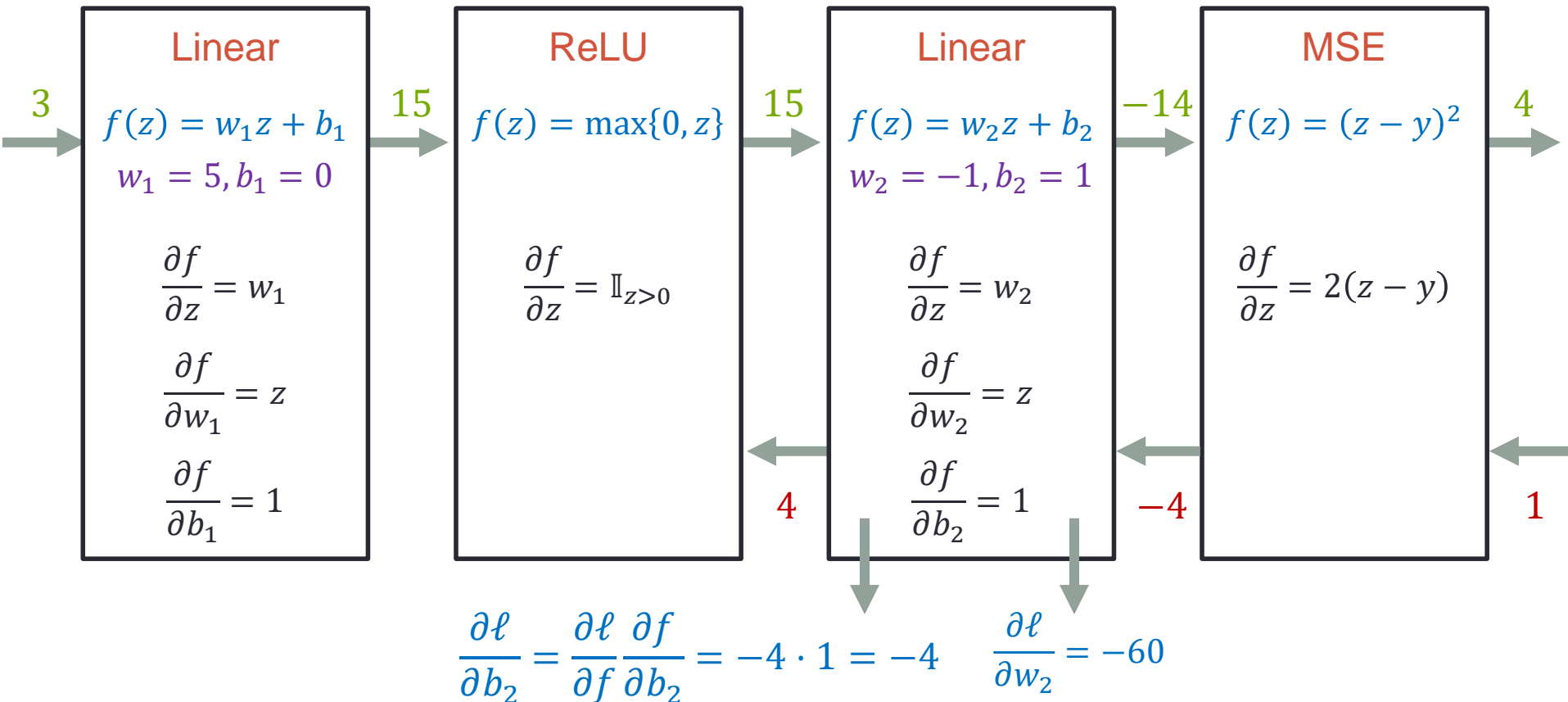


$$\frac{\partial \ell}{\partial z} = \frac{\partial \ell}{\partial f} \frac{\partial f}{\partial z} = -4 \cdot (-1)$$

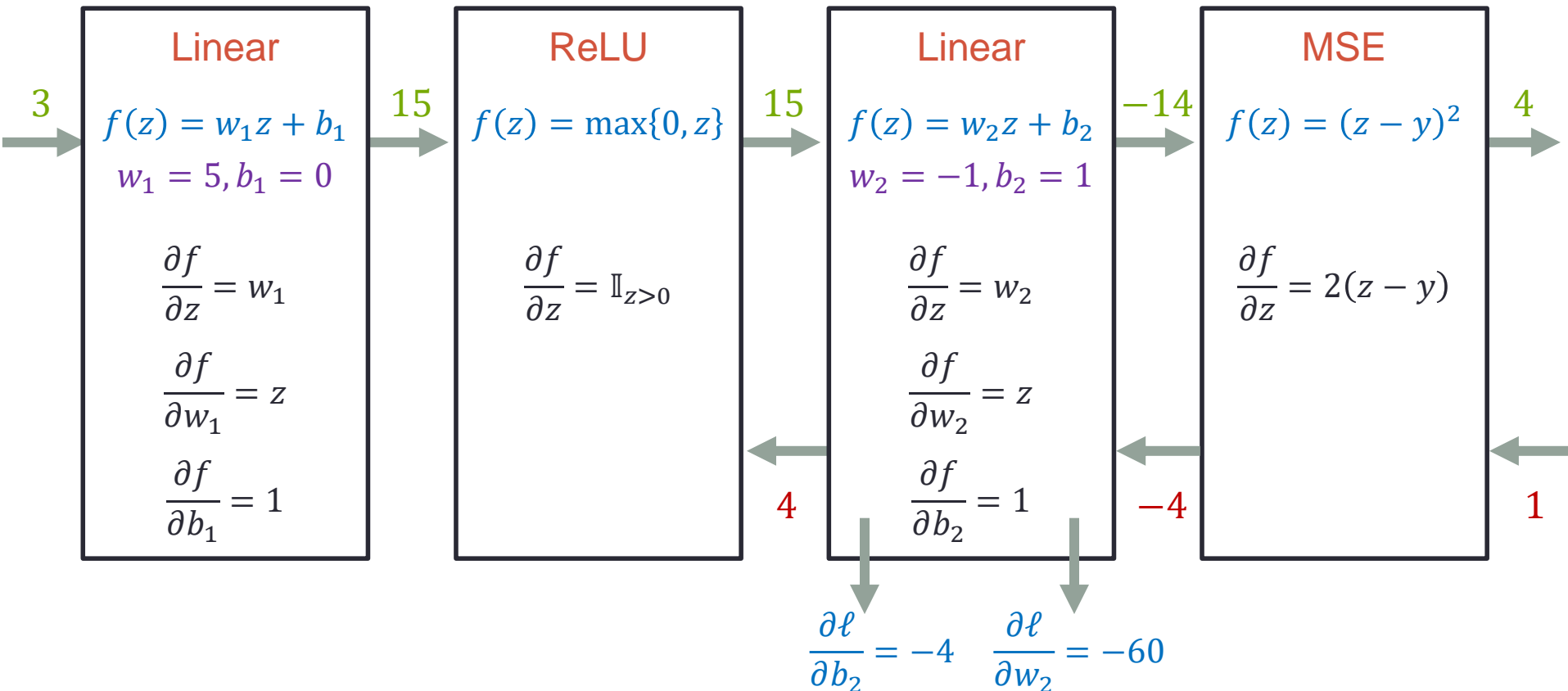
- Initialization: $w_1 = 5, w_2 = -1, b_1 = 0, b_2 = 1$
- Input: $x = 3, y = -12$
- First, we compute the **forward** pass.
- Then, to compute gradients w.r.t θ , we perform a **backward** pass.



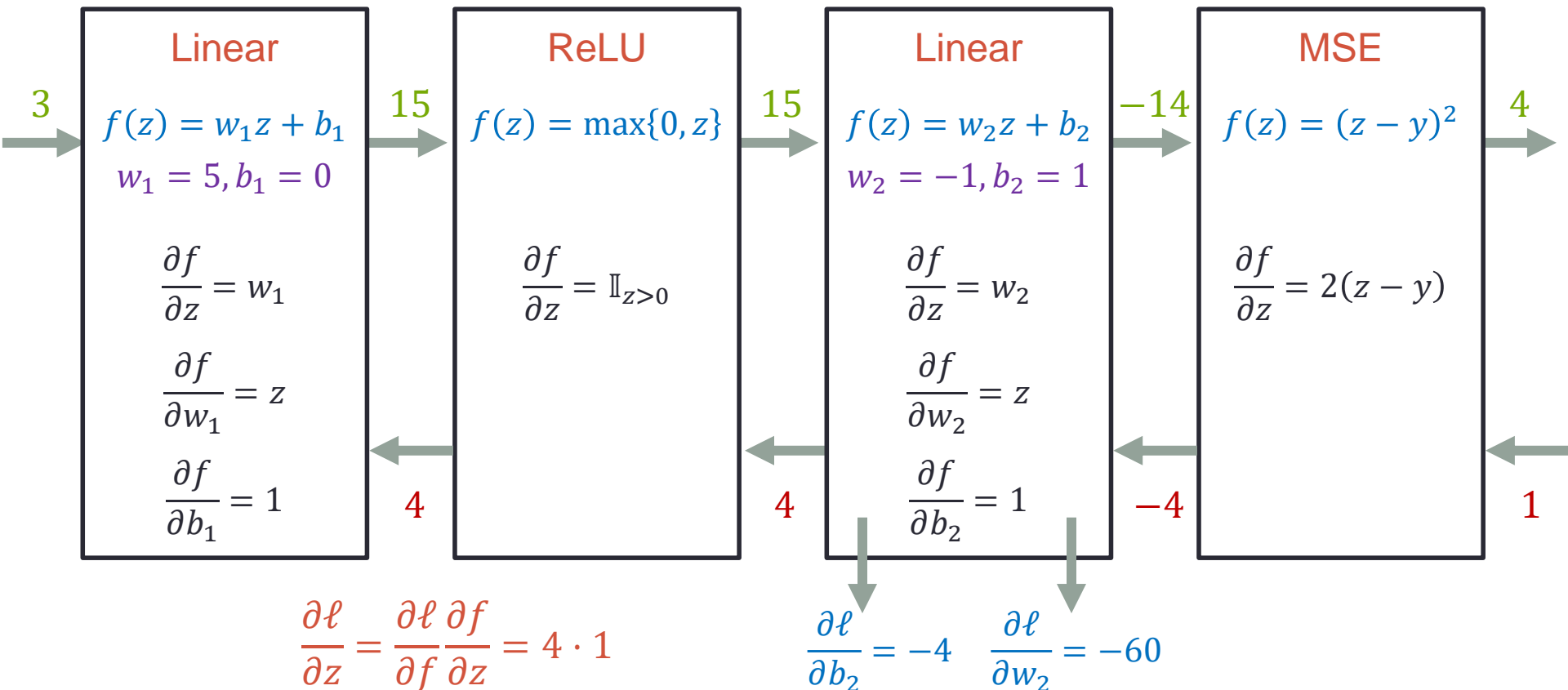
- Initialization: $w_1 = 5, w_2 = -1, b_1 = 0, b_2 = 1$
- Input: $x = 3, y = -12$
- First, we compute the **forward** pass.
- Then, to compute gradients w.r.t θ , we perform a **backward** pass.



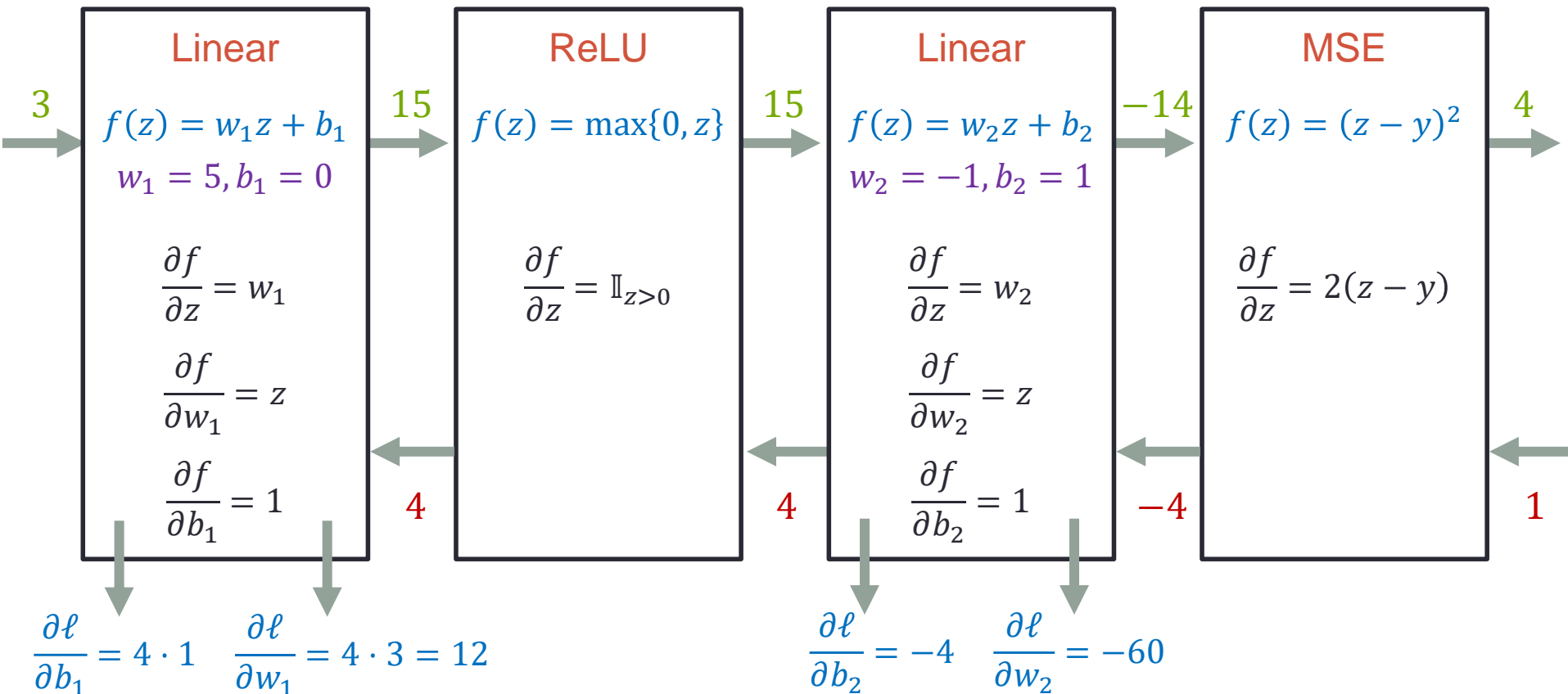
- Initialization: $w_1 = 5, w_2 = -1, b_1 = 0, b_2 = 1$
- Input: $x = 3, y = -12$
- First, we compute the **forward** pass.
- Then, to compute gradients w.r.t θ , we perform a **backward** pass.



- Initialization: $w_1 = 5, w_2 = -1, b_1 = 0, b_2 = 1$
- Input: $x = 3, y = -12$
- First, we compute the **forward** pass.
- Then, to compute gradients w.r.t θ , we perform a **backward** pass.



- Initialization: $w_1 = 5, w_2 = -1, b_1 = 0, b_2 = 1$
- Input: $x = 3, y = -12$
- First, we compute the **forward** pass.
- Then, to compute gradients w.r.t θ , we perform a **backward** pass.



Let's play!

- [Simple regression task](#)
- [Tensorflow playground](#) (extra)

Summary

- Neural networks can represent complicated functions
- Backpropagation:
 - Efficiently computes the gradients of the loss w.r.t the network parameters.
 - Widely used for gradient optimization methods
- Recommended read/watch: [3Blue1Brown – Neural networks](#)