# An Introduction to Deep Learning
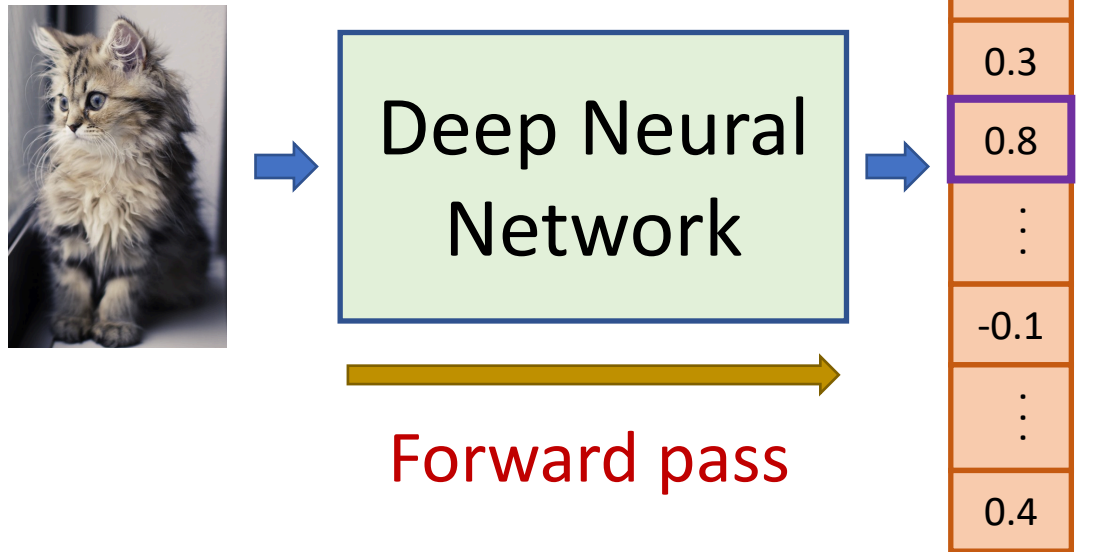
## (part 2)

University of BRISTOL

Farnoosh Heidarivincheh
EMAT31530 - February 2021

# Network Training

So far, we have discussed



**Network output**

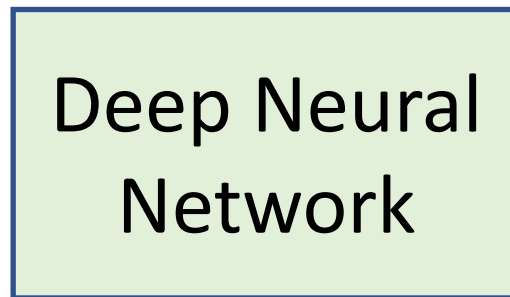| -0.2 |
| 0.3 |
| 0.8 |
| ⋮ |
| -0.1 |
| ⋮ |
| 0.4 |

**Forward pass**

Optimisation during training

1. It makes a prediction using its current parameter set (weights).
2. We give it a feedback of how good its prediction was.
3. It updates (refines) its parameter set according to our feedback

# Network Training

So far, we have discussed



Optimisation during training

1. It makes a prediction using its current parameter set (weights).
2. **We give it a feedback of how good its prediction was.**
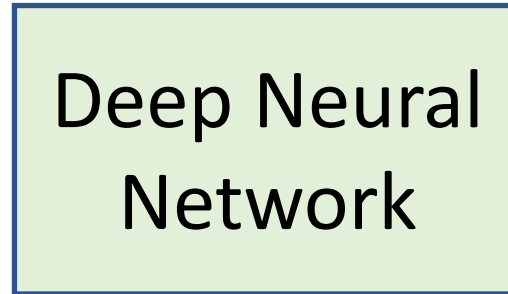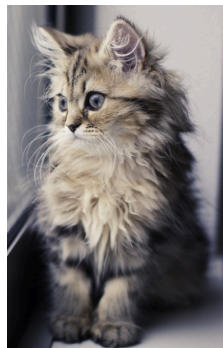3. It updates (refines) its parameter set according to our feedback

# Outline

Network Training

- Gradient Descent

- Back Propagation

# Network Training

# (Gradient Descent)

# Network Training



Network output

Ground truth

| -0.2 | | 0 | dog |
| 0.3 | | 0 | tiger |
| 0.8 | | 0 | ostrich |
| ⋮ | | ⋮ | |
| -0.1 | | 1 | Cat |
| ⋮ | | ⋮ | |
| 0.4 | | 0 | horse |

Our objective is to find an optimal parameter set ($\boldsymbol{W}$) which minimises our loss function

$$\boldsymbol{W}^* = \mathrm{argmin}_{\boldsymbol{W}} \sum_{n=1}^{N} L(f(\boldsymbol{x}_n; \boldsymbol{W}), \boldsymbol{g})$$

Optimisation during training

1. It makes a prediction using its current parameter set (weights).
2. We give it a feedback of how good its prediction was.
3. It updates (refines) its parameter set according to our feedback

# Network Training



Network output

Ground truth

| | |
|---|---|
| -0.2 | 0 dog |
| 0.3 | 0 tiger |
| 0.8 | 0 ostrich |
| ⋮ | ⋮ |
| -0.1 | 1 Cat |
| ⋮ | ⋮ |
| 0.4 | 0 horse |

Our objective is to find an optimal parameter set ($\boldsymbol{W}$) which minimises our loss function

$$\boldsymbol{W}^* = \mathrm{argmin}_{\boldsymbol{W}} \sum_{n=1}^{N} L(f(\boldsymbol{x}_n; \boldsymbol{W}), \boldsymbol{g})$$

Optimisation during training

1. It makes a prediction using its current parameter set (weights).
2. We give it a feedback of how good its prediction was.
3. It updates (refines) its parameter set according to our feedback

# Network Training



Network output

Ground truth

| | |
|---|---|
| -0.2 | |
| 0.3 | |
| 0.8 | |
| ⋮ | |
| -0.1 | |
| ⋮ | |
| 0.4 | |

| | |
|---|---|
| 0 | dog |
| 0 | tiger |
| 0 | ostrich |
| ⋮ | |
| 1 | Cat |
| ⋮ | |
| 0 | horse |

Our objective is to find an optimal parameter set ($\boldsymbol{W}$) which minimises our loss function

$$\boldsymbol{W}^* = \mathrm{argmin}_{\boldsymbol{W}} \sum_{n=1}^{N} L(f(\boldsymbol{x}_n; \boldsymbol{W}), \boldsymbol{g})$$
$$= \mathrm{argmin}_{\boldsymbol{W}} J(\boldsymbol{W})$$
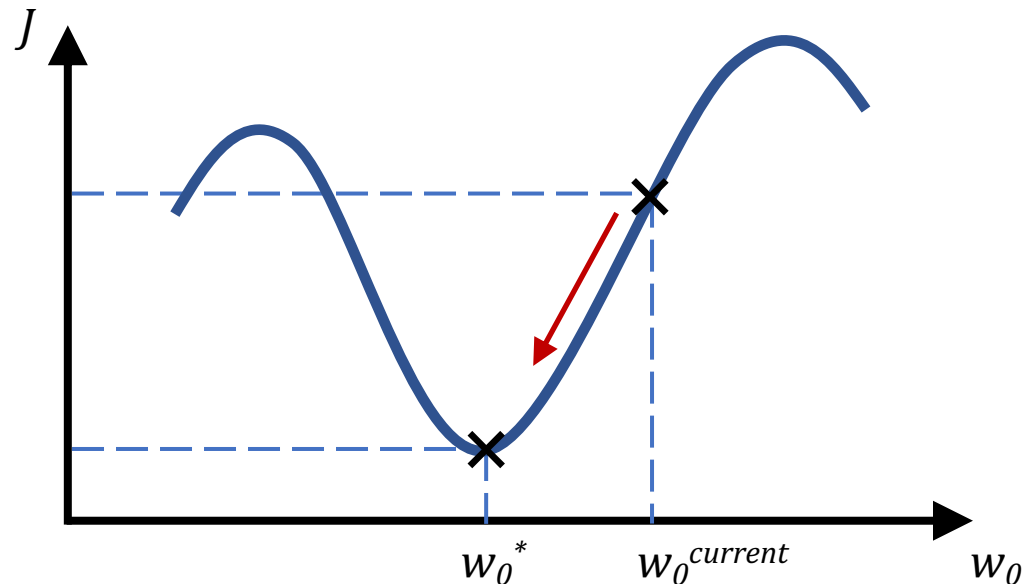
Optimisation during training

1. It makes a prediction using its current parameter set (weights).
2. We give it a feedback of how good its prediction was.
3. It updates (refines) its parameter set according to our feedback

# Network Training

- We start from our current $\boldsymbol{W}$

- We need to move towards the direction which minimises the loss

$$\boldsymbol{W}^* = \mathrm{argmin}_{\boldsymbol{W}} \sum_{n=1}^{N} L(f(\boldsymbol{x}_n; \boldsymbol{W}), \boldsymbol{g})$$

$$= \mathrm{argmin}_{\boldsymbol{W}} J(\boldsymbol{W})$$

Note: $\boldsymbol{W}$ is the aggregation of several multi-dimensional matrices

# Network Training

- We start from our current $W$

- We need to move towards the direction which minimises the loss

$\Longrightarrow$ **Gradient Descent (GD)**

We take the gradient of the loss w.r.t the current $W$
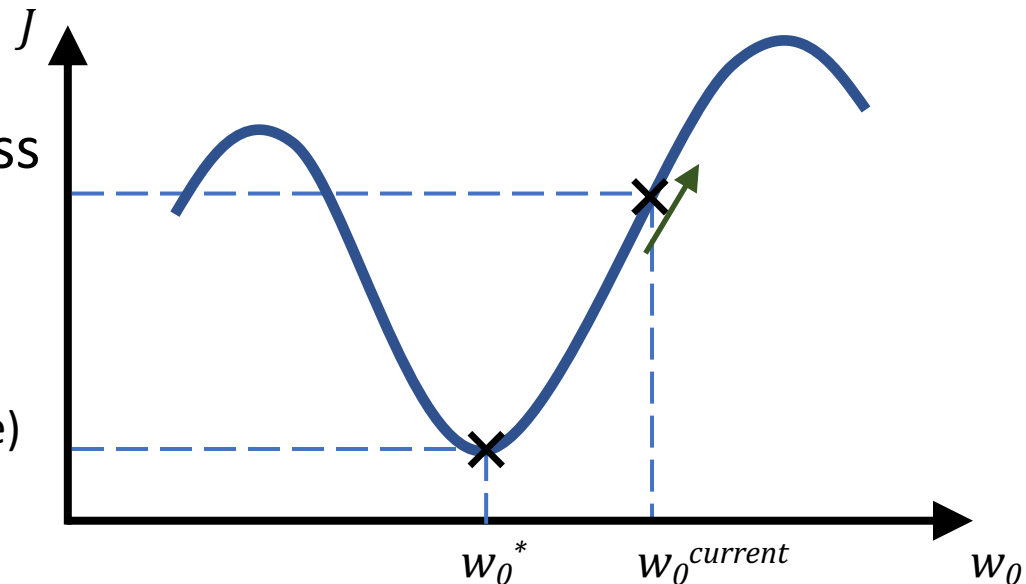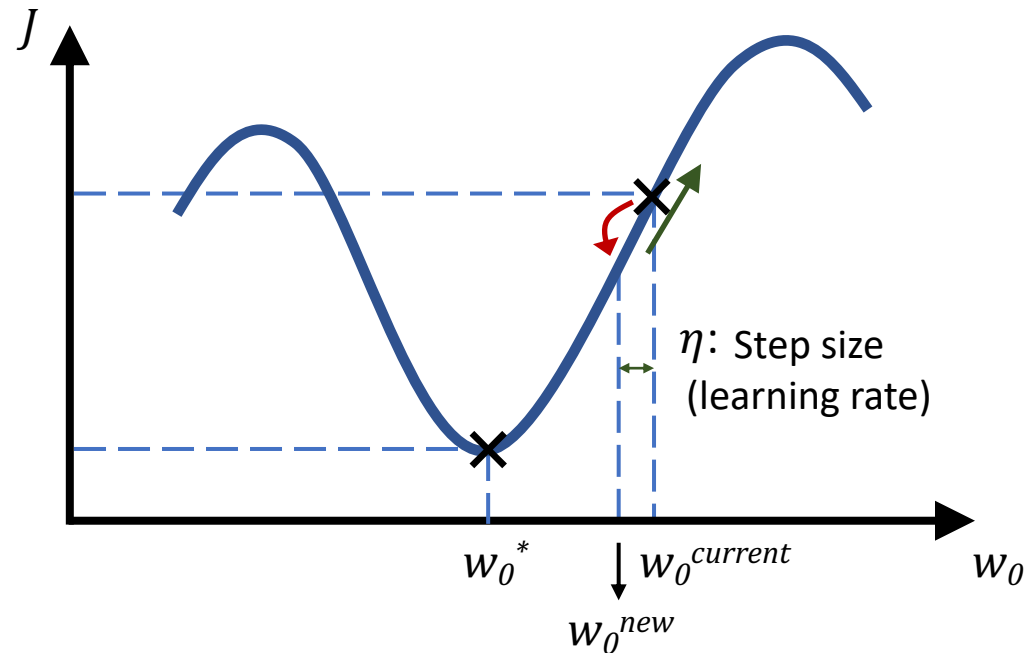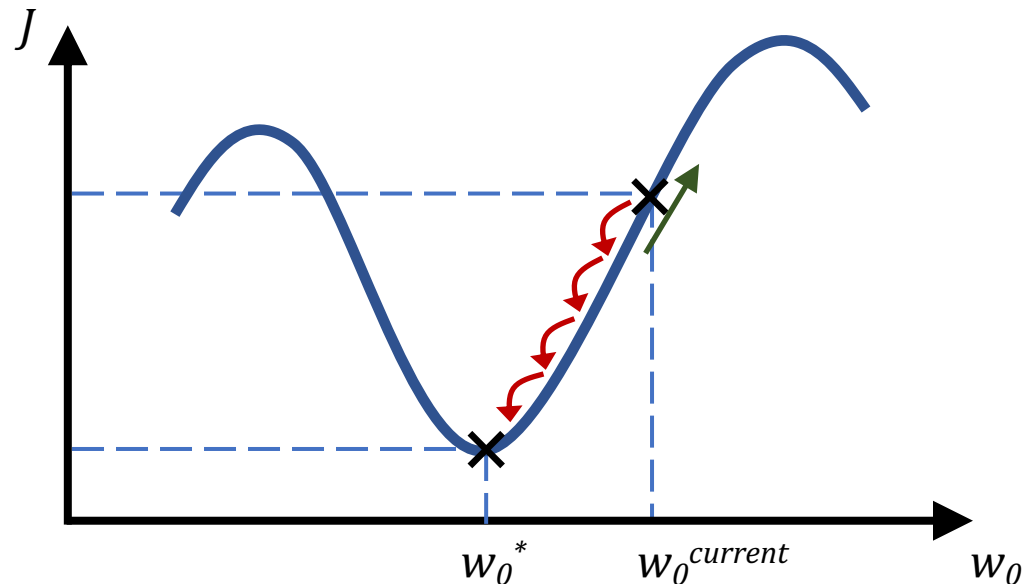
(The loss function must be differentiable)

# Network Training

- We start from our current $W$

- We need to move towards the direction which minimises the loss

⟹ **Gradient Descent (GD)**

We step in the opposite direction of the gradient

$$W^{new} = W^{current} - \eta \frac{\partial J}{\partial W}$$



$\eta$: Step size (learning rate)

$w_0^*$     $\downarrow w_0^{current}$     $w_0$

$w_0^{new}$

# Network Training

- We start from our current $W$

- We need to move towards the direction which minimises the loss

⇒ **Gradient Descent (GD)**

We repeat until convergence

# Network Training

**Stochastic Gradient Descent (SGD)**

- GD is expensive to be computed on the whole training data ($\boldsymbol{X}$)

- We can use a single data point, e.g. one image

- This data point is selected randomly

$$\boldsymbol{W}^{new} = \boldsymbol{W}^{current} - \eta \frac{\partial J(\boldsymbol{W}; \boldsymbol{x}_i)}{\partial \boldsymbol{W}}$$

where, $\quad \boldsymbol{x}_i \in \boldsymbol{X}$

# Network Training

**Mini-Batch Gradient Descent**

- We divide our training set into mini-batches

- Each mini-batch contains $B$ data points

$$\boldsymbol{W}^{new} = \boldsymbol{W}^{current} - \eta \frac{\partial J(\boldsymbol{W}; \boldsymbol{x}_{Batch})}{\partial \boldsymbol{W}}$$

where, $\boldsymbol{x}_{Batch} \subset \boldsymbol{X}$ and $|\boldsymbol{x}_{Batch}| = B$

↓

(Batch size)

# Network Training

# (Back Propagation)

# Network Training

- $W$ is the aggregation of several multi-dimensional matrices $W^{(1)}$, $W^{(2)}$, ..., $W^{(n)}$

- The network input goes through *layers* one after another

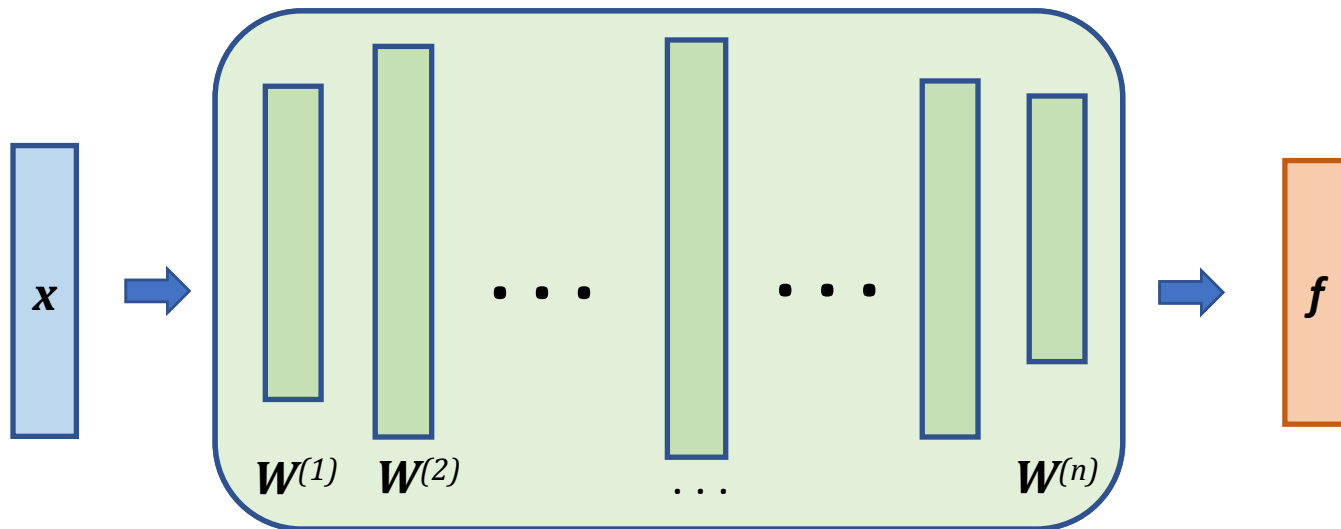How should we compute the loss gradient w.r.t. $W$ in all those layers?

# Network Training

- $W$ is the aggregation of several multi-dimensional matrices $W^{(1)}$, $W^{(2)}$, …, $W^{(n)}$

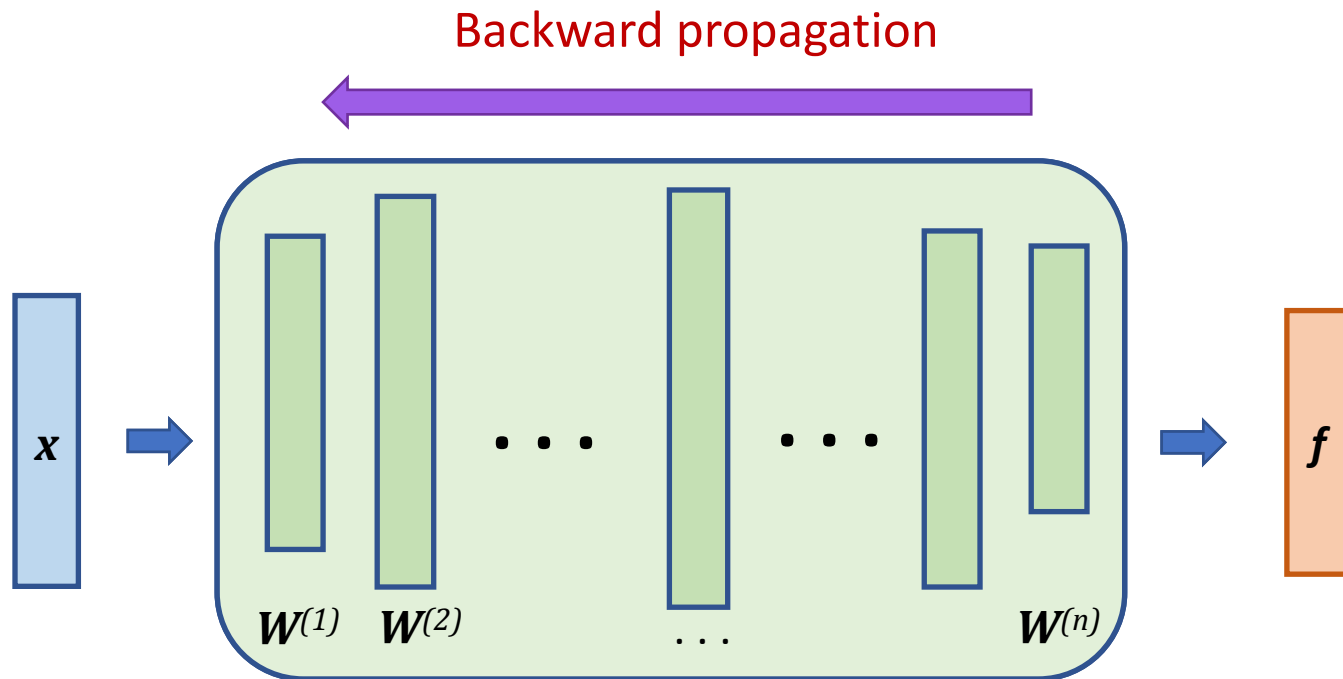- The network input goes through *layers* one after another

How should we compute the loss gradient w.r.t. $W$ in all those layers?
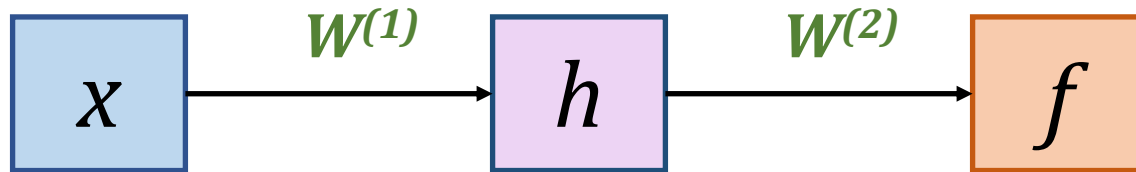
**Back Propagation**

# Network Training

Back Propagation (BP) uses **chain rule** recursively to back-propagate our gradient calculations **layer by layer** from **right to left**

Backward propagation

$x$   →   ... ... →   $f$

$W^{(1)}$    $W^{(2)}$      ...      $W^{(n)}$

# Network Training

Example: 2 neurons with 1D input



$W^{(1)} = [w_1, w_{b1}] = [3, -1]$

$W^{(2)} = [w_2, w_{b2}] = [2, 4]$

$x = -2$

We want to calculate $\dfrac{\partial f}{\partial w_1}, \dfrac{\partial f}{\partial w_{b1}}, \dfrac{\partial f}{\partial w_2}$ and $\dfrac{\partial f}{\partial w_{b2}}$
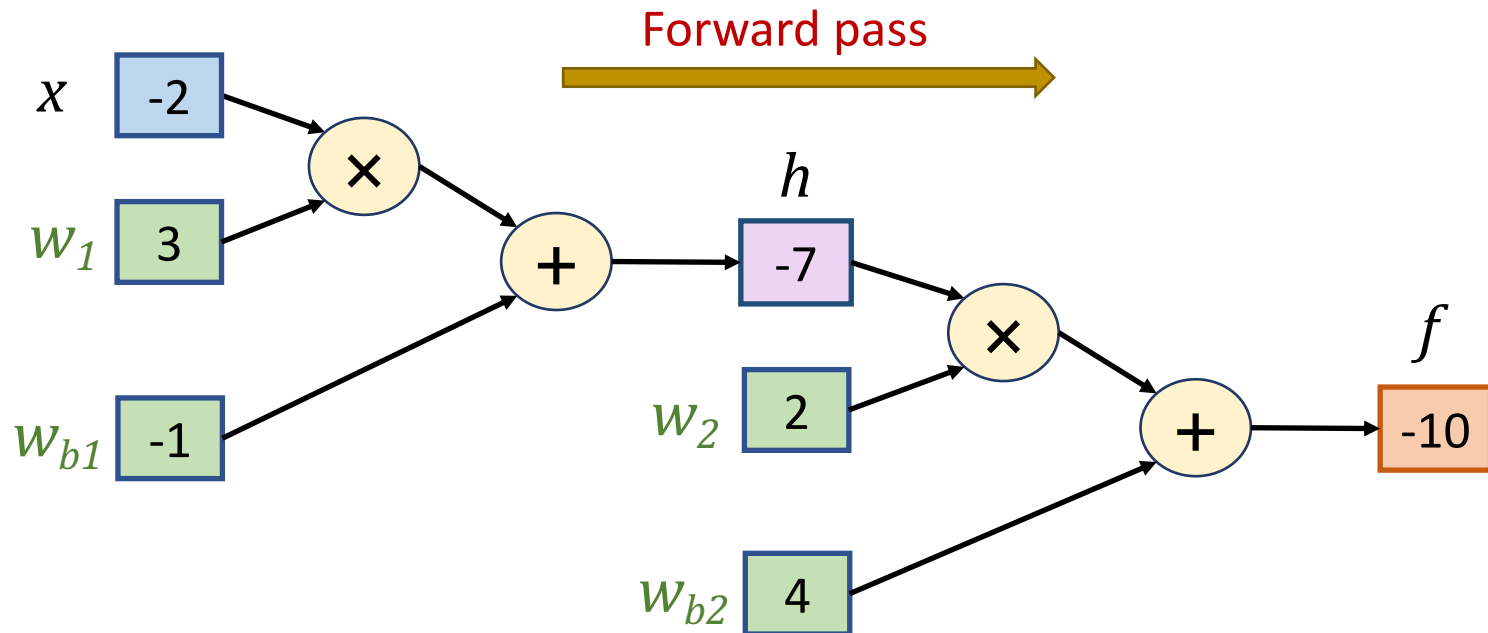
# Network Training

Example: 2 neurons with 1D input

$$h = xw_1 + w_{b1}$$

$$f = hw_2 + w_{b2}$$

(For simplicity, we remove the activation and loss functions)

Forward pass

$x$ [-2]

$w_1$ [3]

$w_{b1}$ [-1]

$h$ [-7]
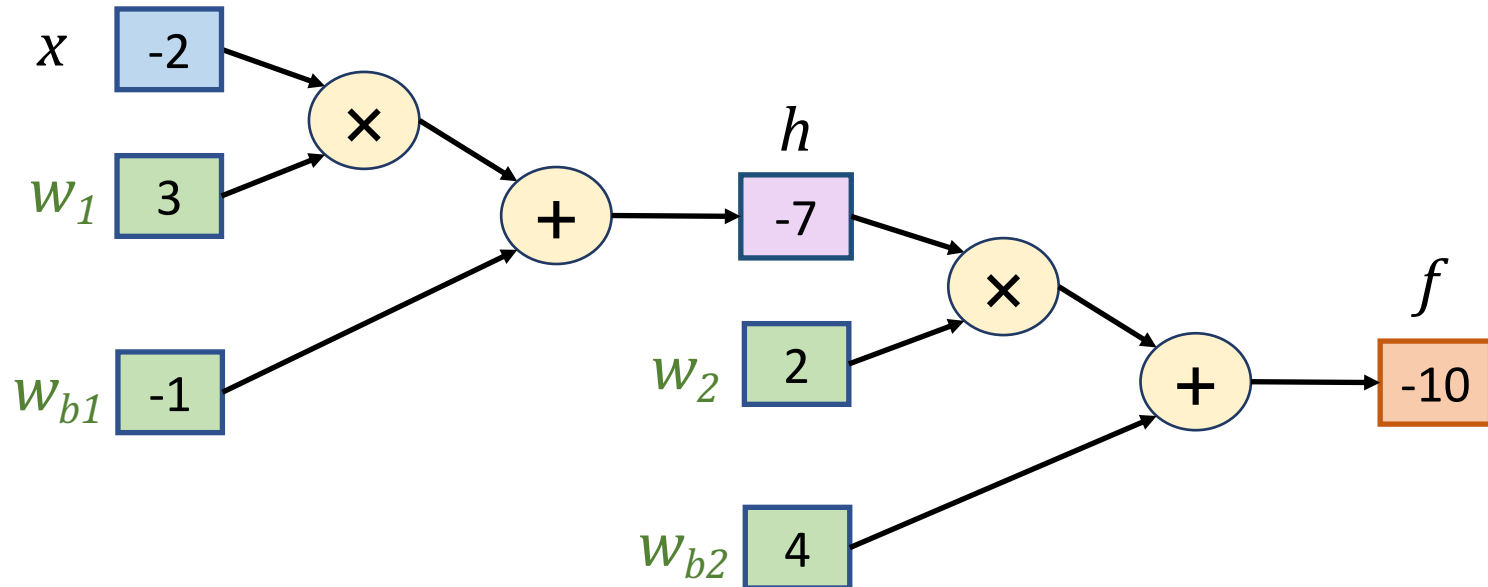
$w_2$ [2]

$w_{b2}$ [4]

$f$ [-10]

# Network Training

Example: 2 neurons with 1D input

$$h = xw_1 + w_{b1} = -2{\times}3 - 1 = -7$$

$$f = hw_2 + w_{b2} = -7{\times}2 + 4 = -10$$

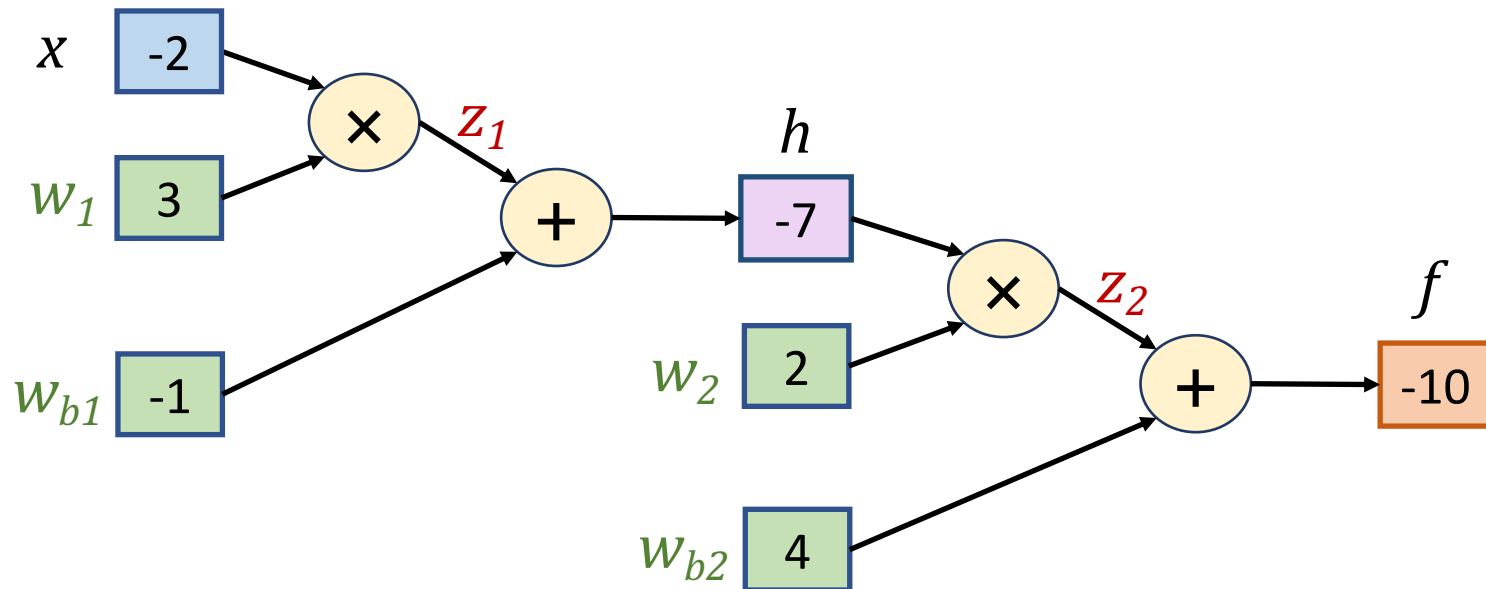# Network Training

Example: 2 neurons with 1D input

$$h = \underset{z_1}{\boxed{(xw_1)}} + w_{b1} = -2 \times 3 - 1 = -7$$

$$f = \underset{z_2}{\boxed{(hw_2)}} + w_{b2} = -7 \times 2 + 4 = -10$$

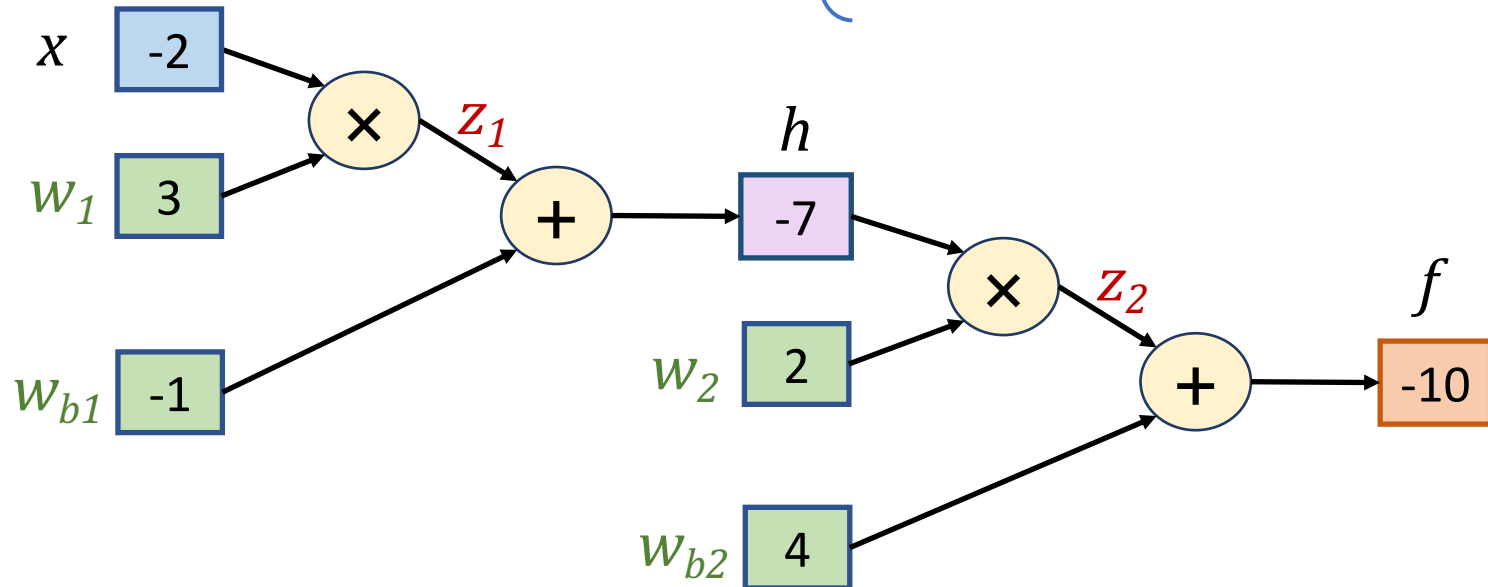# Network Training

Example: 2 neurons with 1D input

$$h = xw_1 + w_{b1} = -2 \times 3 - 1 = -7$$

$$f = hw_2 + w_{b2} = -7 \times 2 + 4 = -10$$

$$\begin{cases} \dfrac{\partial f}{\partial w_{b2}} = 1 \\ \\ \dfrac{\partial f}{\partial w_2} = \dfrac{\partial f}{\partial z_2} \times \dfrac{\partial z_2}{\partial w_2} = 1 \times h = -7 \end{cases}$$

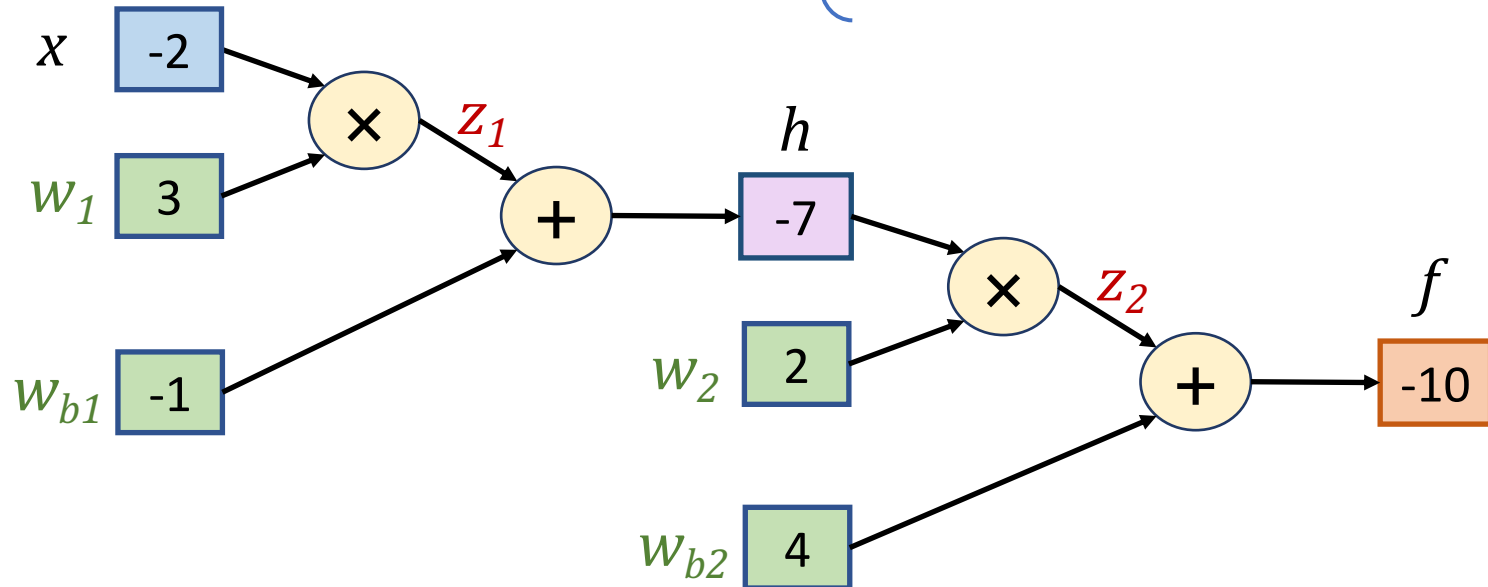# Network Training

Example: 2 neurons with 1D input

$$h = xw_1 + w_{b1} = -2 \times 3 - 1 = -7$$

$$f = hw_2 + w_{b2} = -7 \times 2 + 4 = -10$$

$$\frac{\partial f}{\partial w_{b2}} = 1$$

$$\frac{\partial f}{\partial w_2} = \frac{\partial f}{\partial z_2} \times \frac{\partial z_2}{\partial w_2} = 1 \times h = -7$$

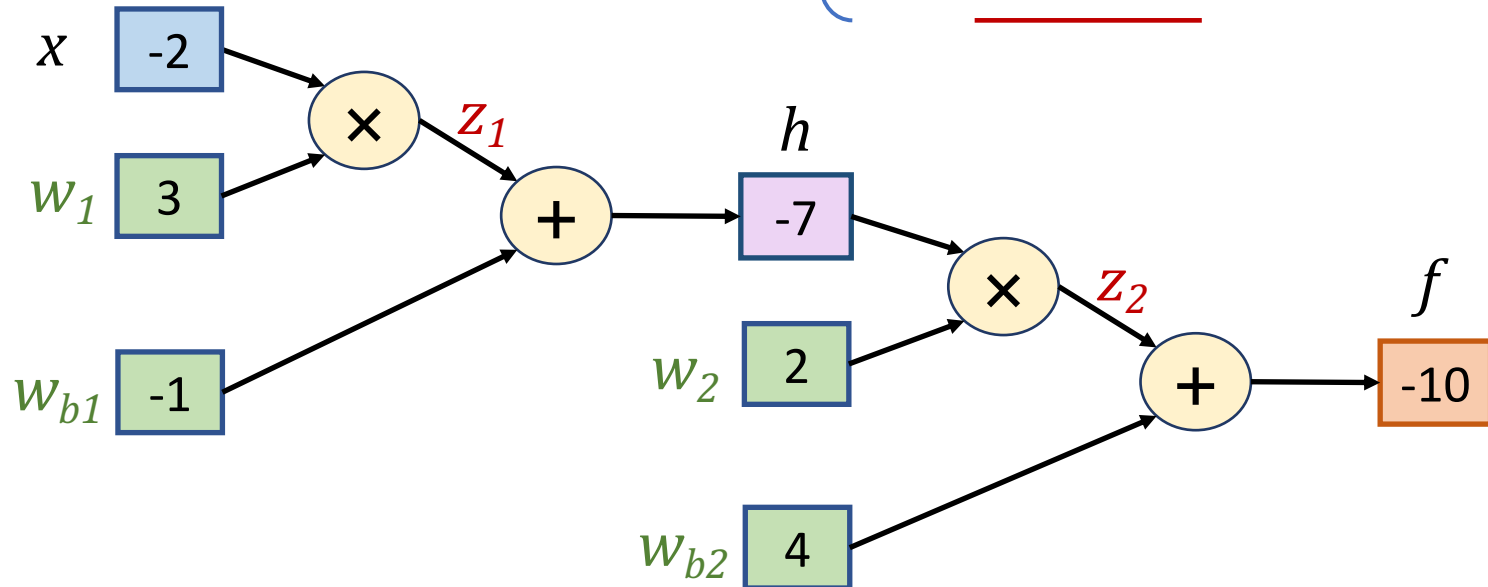# Network Training

Example: 2 neurons with 1D input

$$h = xw_1 + w_{b1} = -2 \times 3 - 1 = -7$$

$$f = hw_2 + w_{b2} = -7 \times 2 + 4 = -10$$

$$\underline{z_2}$$

$$\frac{\partial f}{\partial w_{b2}} = 1$$

$$\frac{\partial f}{\partial w_2} = \frac{\partial f}{\partial z_2} \times \frac{\partial z_2}{\partial w_2} = 1 \times h = -7$$

# Network Training

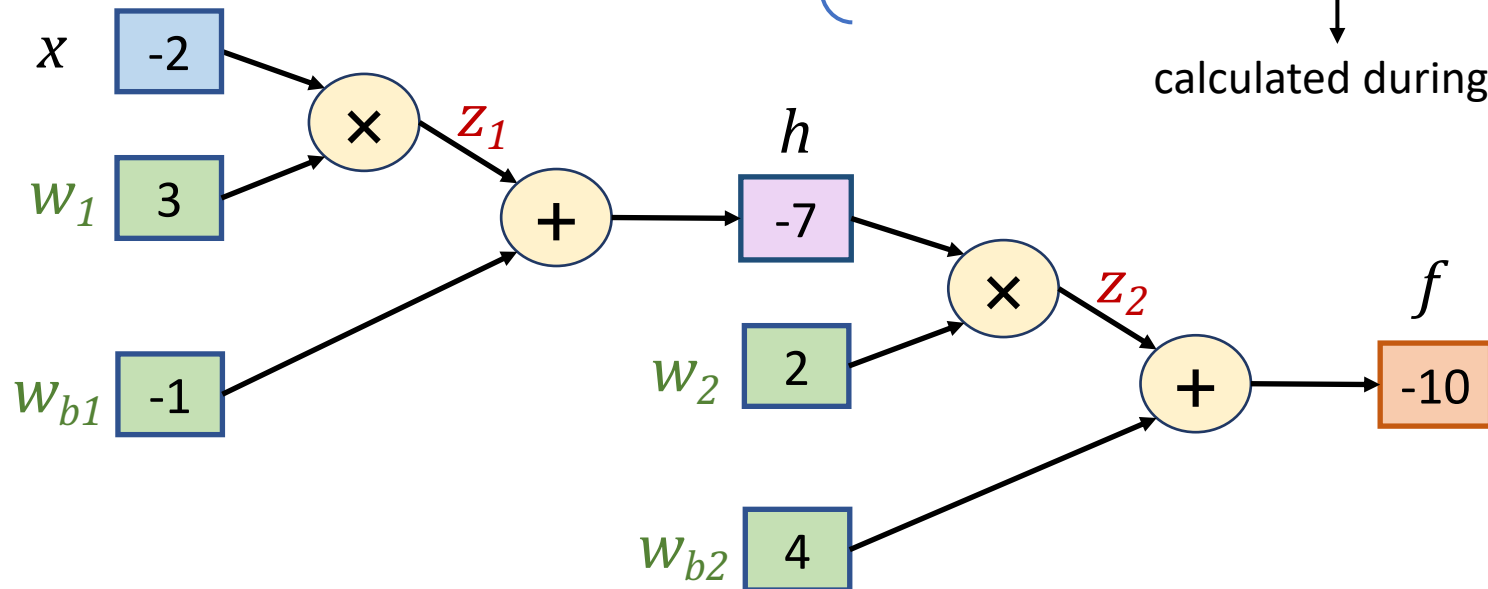Example: 2 neurons with 1D input

$$h = xw_1 + w_{b1} = -2 \times 3 - 1 = -7$$

$$f = hw_2 + w_{b2} = -7 \times 2 + 4 = -10$$

$$\frac{\partial f}{\partial w_{b2}} = 1$$

$$\frac{\partial f}{\partial w_2} = \frac{\partial f}{\partial z_2} \times \frac{\partial z_2}{\partial w_2} = 1 \times h = -7$$
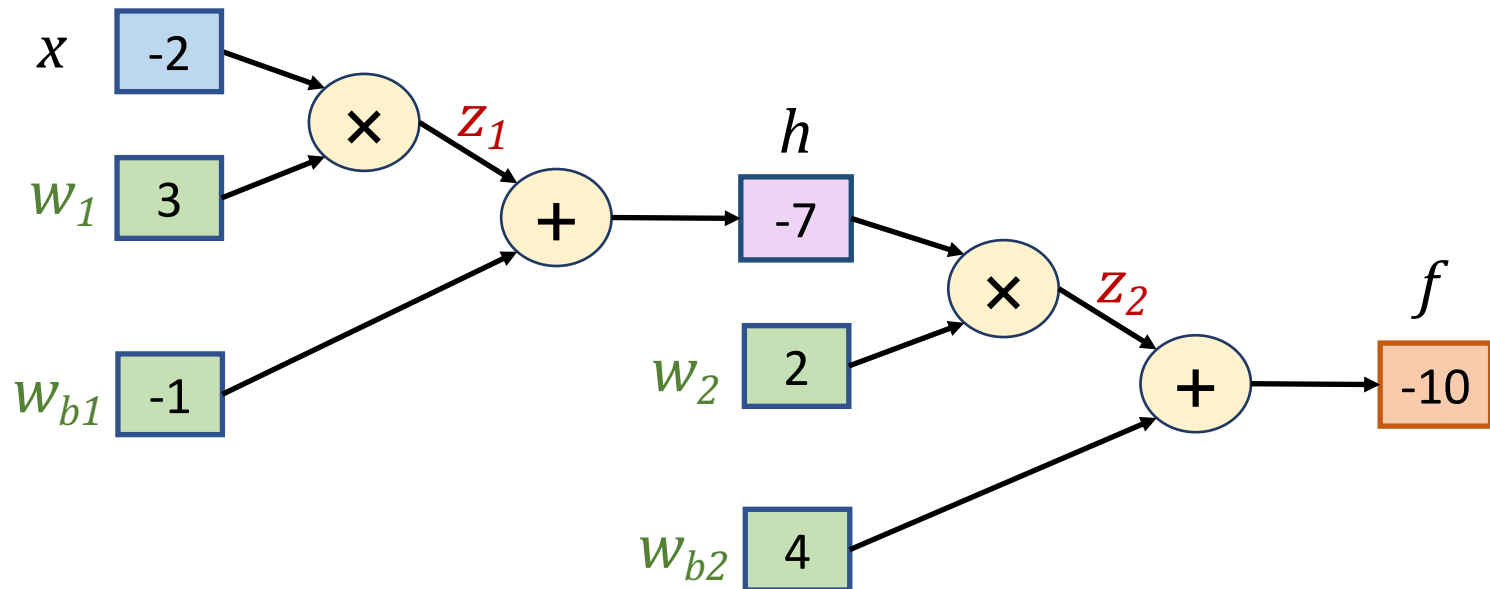
calculated during FP

# Network Training

Example: 2 neurons with 1D input

$$h = xw_1 + w_{b1} = -2 \times 3 - 1 = -7$$

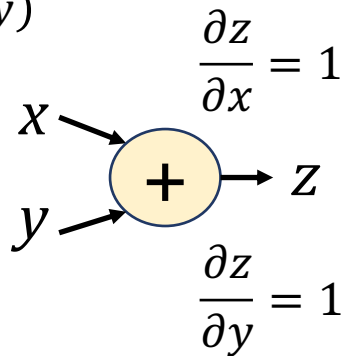$$f = hw_2 + w_{b2} = -7 \times 2 + 4 = -10$$

$$\frac{\partial f}{\partial w_{b1}} = \frac{\partial f}{\partial z_2} \times \frac{\partial z_2}{\partial h} \times \frac{\partial h}{\partial w_{b1}}$$
$$= 1 \times w_2 \times 1 = 2$$

$$\frac{\partial f}{\partial w_1} = \frac{\partial f}{\partial z_2} \times \frac{\partial z_2}{\partial h} \times \frac{\partial h}{\partial z_1} \times \frac{\partial z_1}{\partial w_1}$$
$$= 1 \times w_2 \times 1 \times x = -4$$

# Network Training

- BP computes the partial derivatives of each function **locally** w.r.t. the inputs
- Partial derivates from the previous layers are also multiplied using chain rule

$(z = x + y)$

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

$x$

$+$ $\rightarrow$ $z$

$y$

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

$(z = xy)$

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

$x$

$\times$ $\rightarrow$ $z$
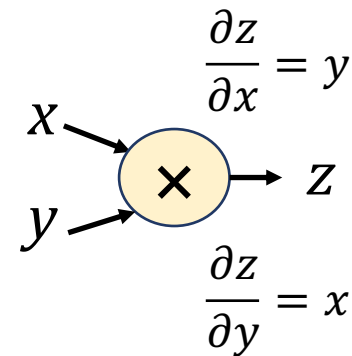
$y$

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

# Network Training

- BP computes the partial derivatives of each function **locally** w.r.t. the inputs
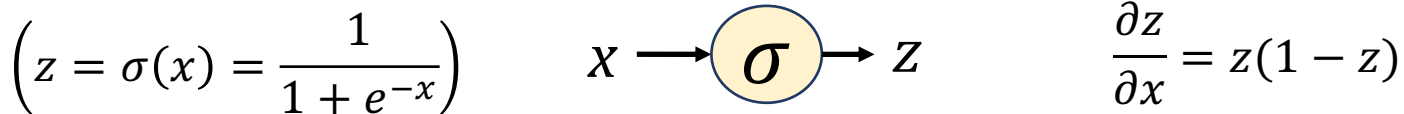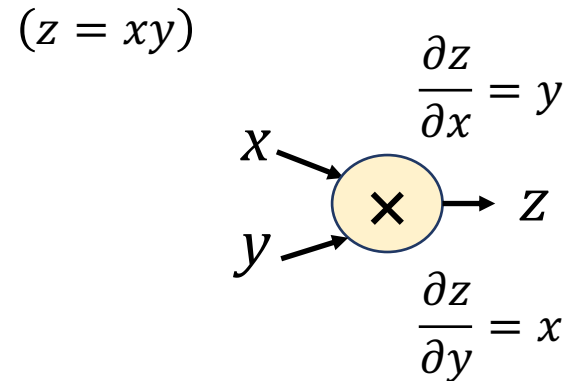- Partial derivates from the previous layers are also multiplied using chain rule

$(z = x + y)$

$x$

$+$

$y$

$z$

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

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

$(z = xy)$

$x$

$\times$

$y$

$z$

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

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

$\left( z = \sigma(x) = \frac{1}{1 + e^{-x}} \right)$

$x$

$\sigma$

$z$

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

# Network Training

- BP computes the partial derivatives of each function **locally** w.r.t. the inputs
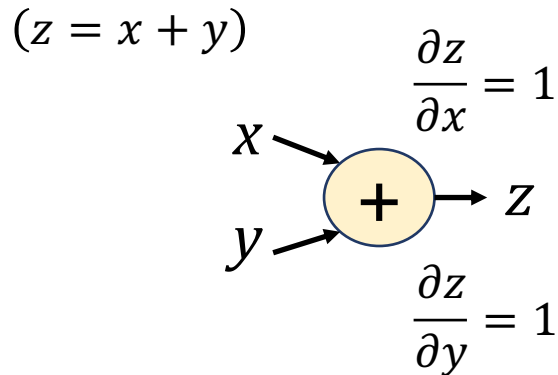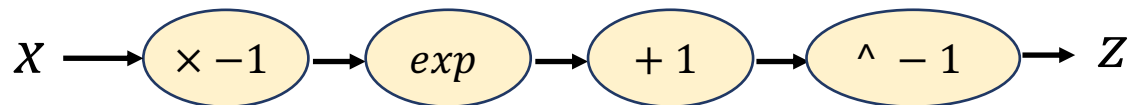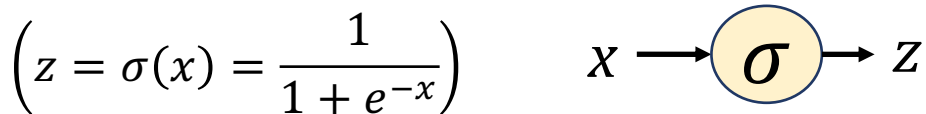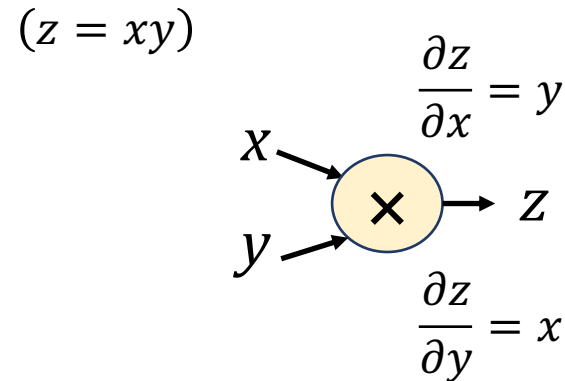- Partial derivates from the previous layers are also multiplied using chain rule

$(z = x + y)$

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

$x \searrow$
$\boxed{+} \rightarrow z$
$y \nearrow$

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

$(z = xy)$

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

$x \searrow$
$\boxed{\times} \rightarrow z$
$y \nearrow$

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

$$\left( z = \sigma(x) = \frac{1}{1 + e^{-x}} \right)$$

$x \rightarrow \boxed{\sigma} \rightarrow z$

$x \rightarrow \boxed{\times -1} \rightarrow \boxed{exp} \rightarrow \boxed{+1} \rightarrow \boxed{\wedge -1} \rightarrow z$
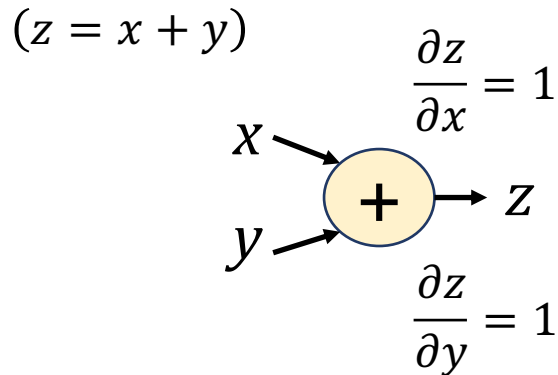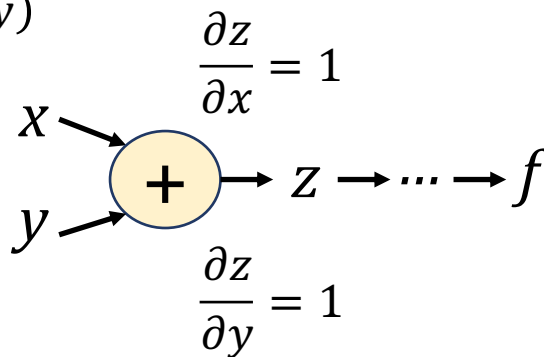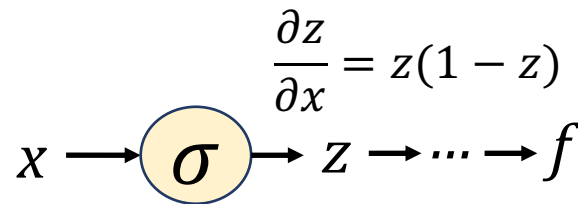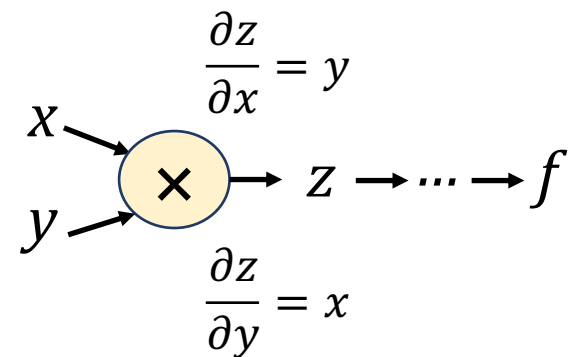
# Network Training

- BP computes the partial derivatives of each function **locally** w.r.t. the inputs
- Partial derivates from the previous layers are also multiplied using chain rule

$(z = x + y)$

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

$$x \searrow$$
$$\boxed{+} \rightarrow z \rightarrow \cdots \rightarrow f$$
$$y \nearrow$$

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

$(z = xy)$

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

$$x \searrow$$
$$\boxed{\times} \rightarrow z \rightarrow \cdots \rightarrow f$$
$$y \nearrow$$

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

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

$$x \rightarrow \boxed{\sigma} \rightarrow z \rightarrow \cdots \rightarrow f$$
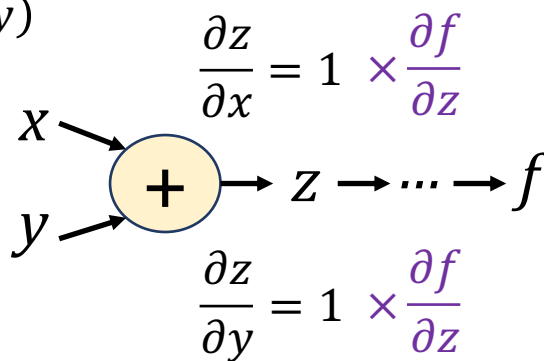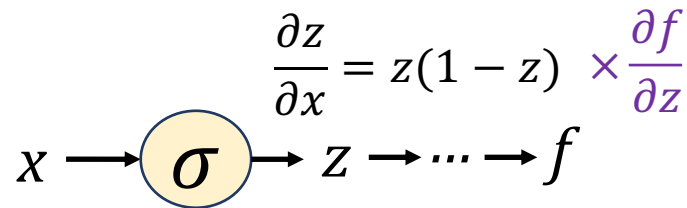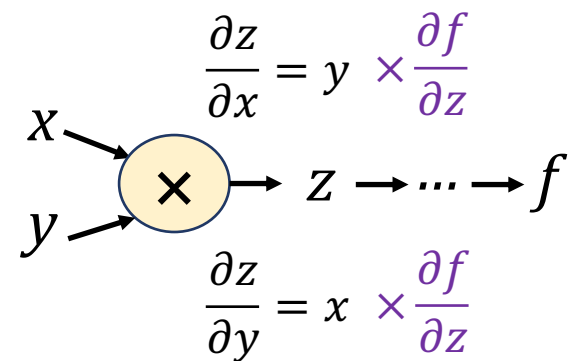
$$\left( z = \sigma(x) = \frac{1}{1 + e^{-x}} \right)$$

# Network Training

- BP computes the partial derivatives of each function **locally** w.r.t. the inputs
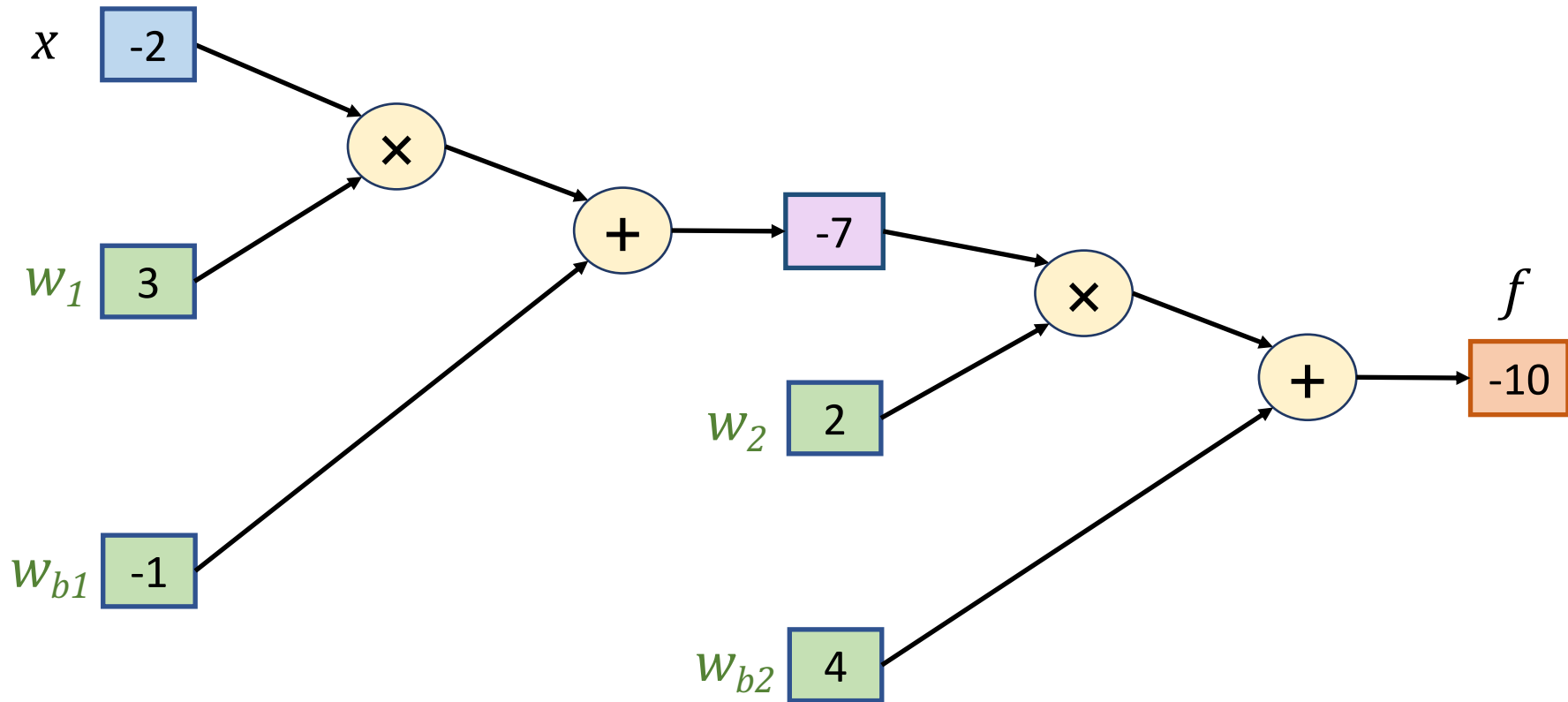- Partial derivates from the previous layers are also multiplied using chain rule
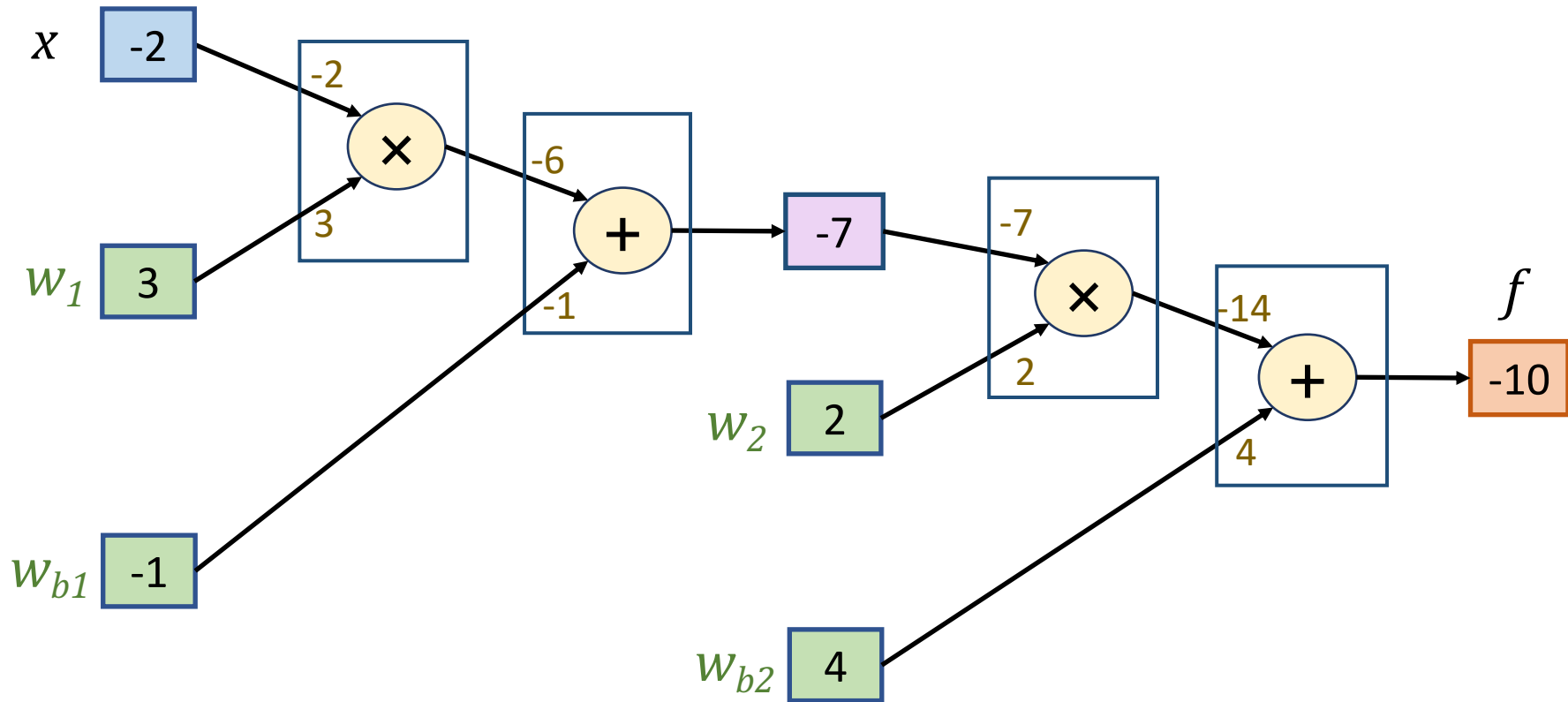
$(z = x + y)$

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

$x \searrow$
$\boxed{+} \rightarrow z \rightarrow \cdots \rightarrow f$
$y \nearrow$

$$\frac{\partial z}{\partial y} = 1 \times \frac{\partial f}{\partial z}$$

$(z = xy)$

$$\frac{\partial z}{\partial x} = y \times \frac{\partial f}{\partial z}$$

$x \searrow$
$\boxed{\times} \rightarrow z \rightarrow \cdots \rightarrow f$
$y \nearrow$

$$\frac{\partial z}{\partial y} = x \times \frac{\partial f}{\partial z}$$

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

$x \rightarrow \boxed{\sigma} \rightarrow z \rightarrow \cdots \rightarrow f$

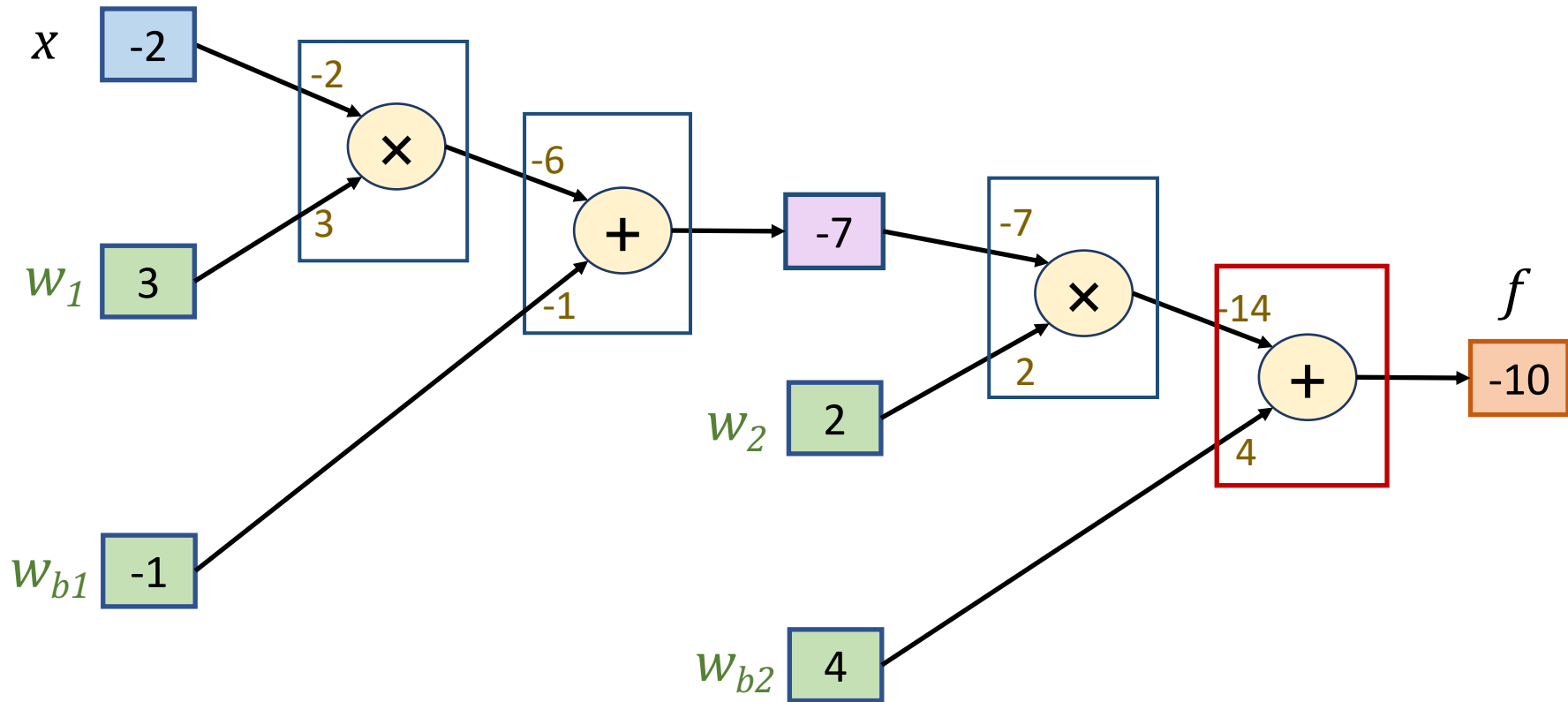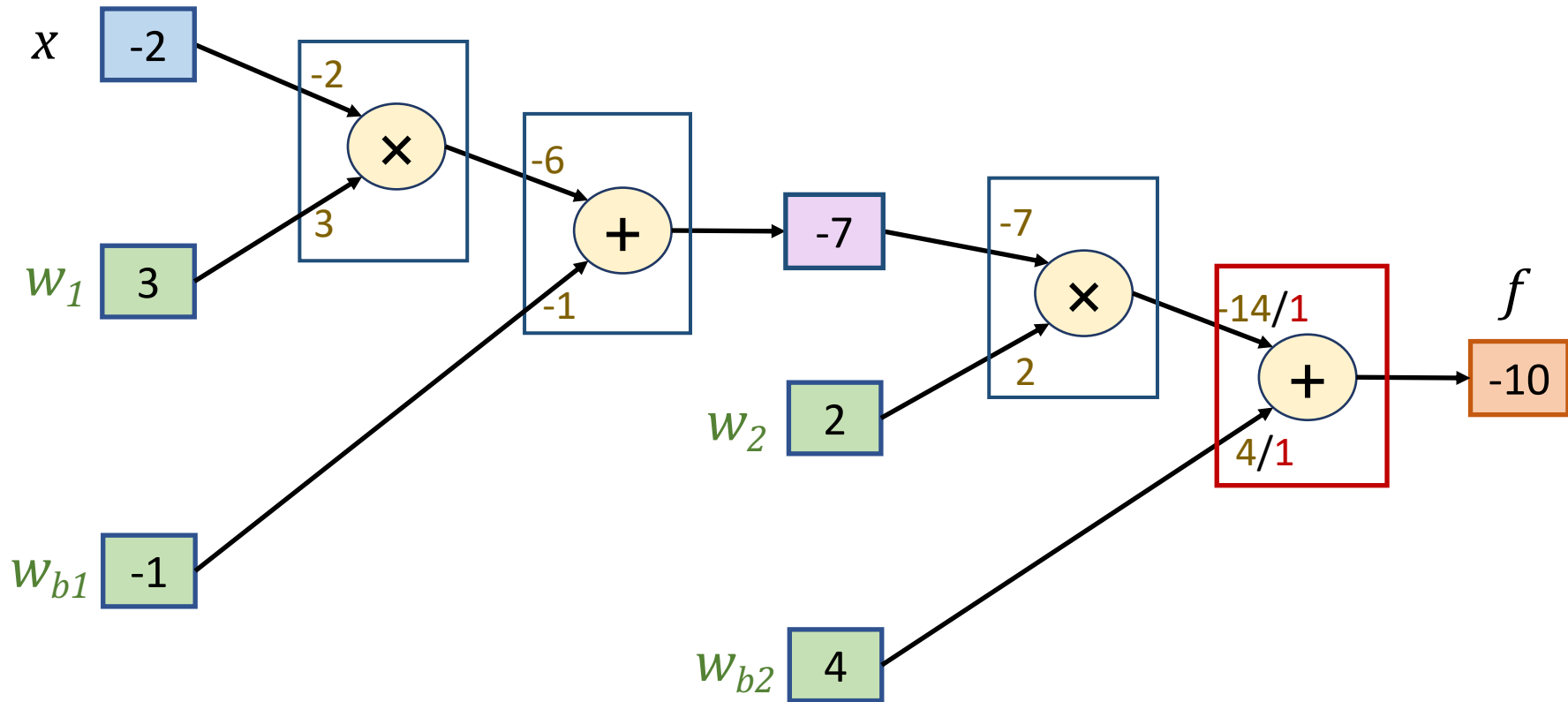$$\left( z = \sigma(x) = \frac{1}{1 + e^{-x}} \right)$$

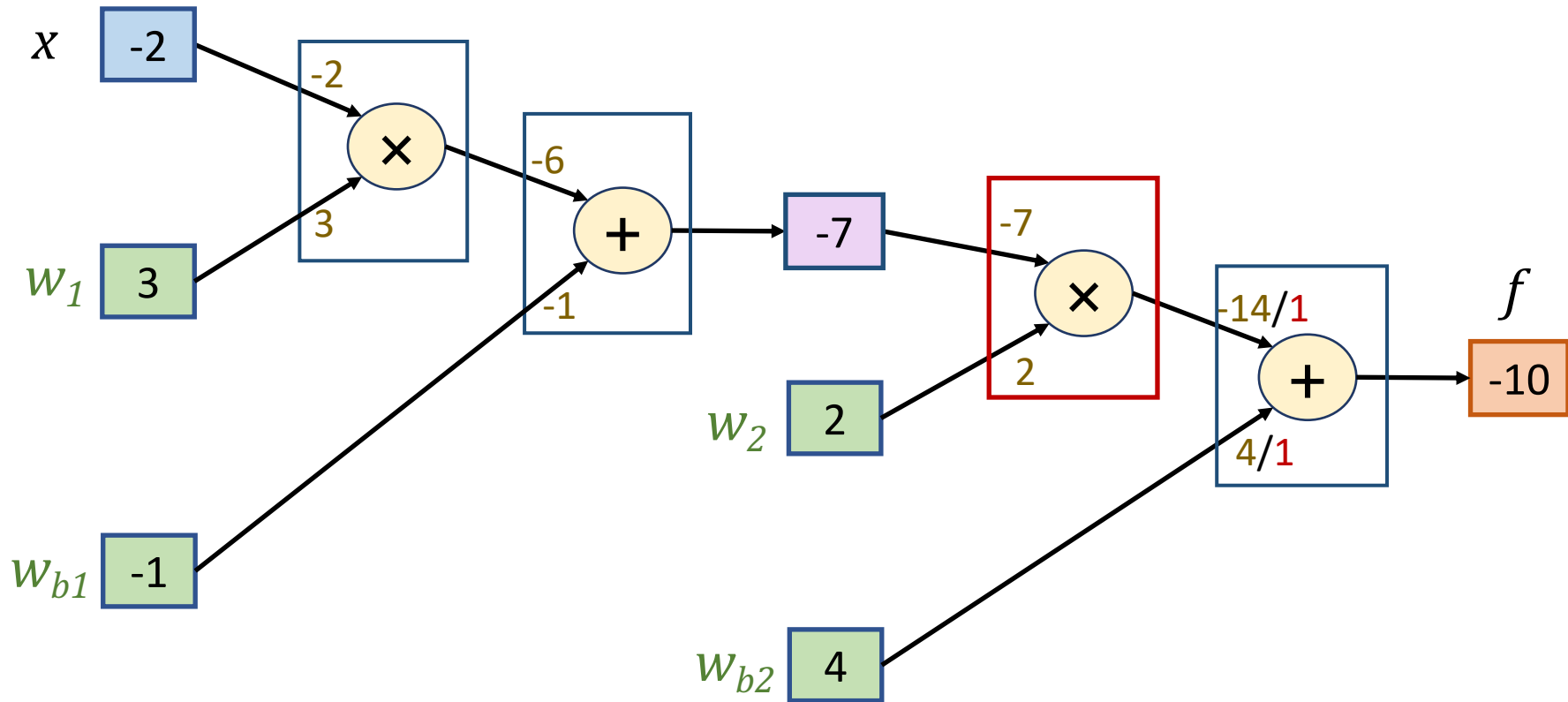# Network Training
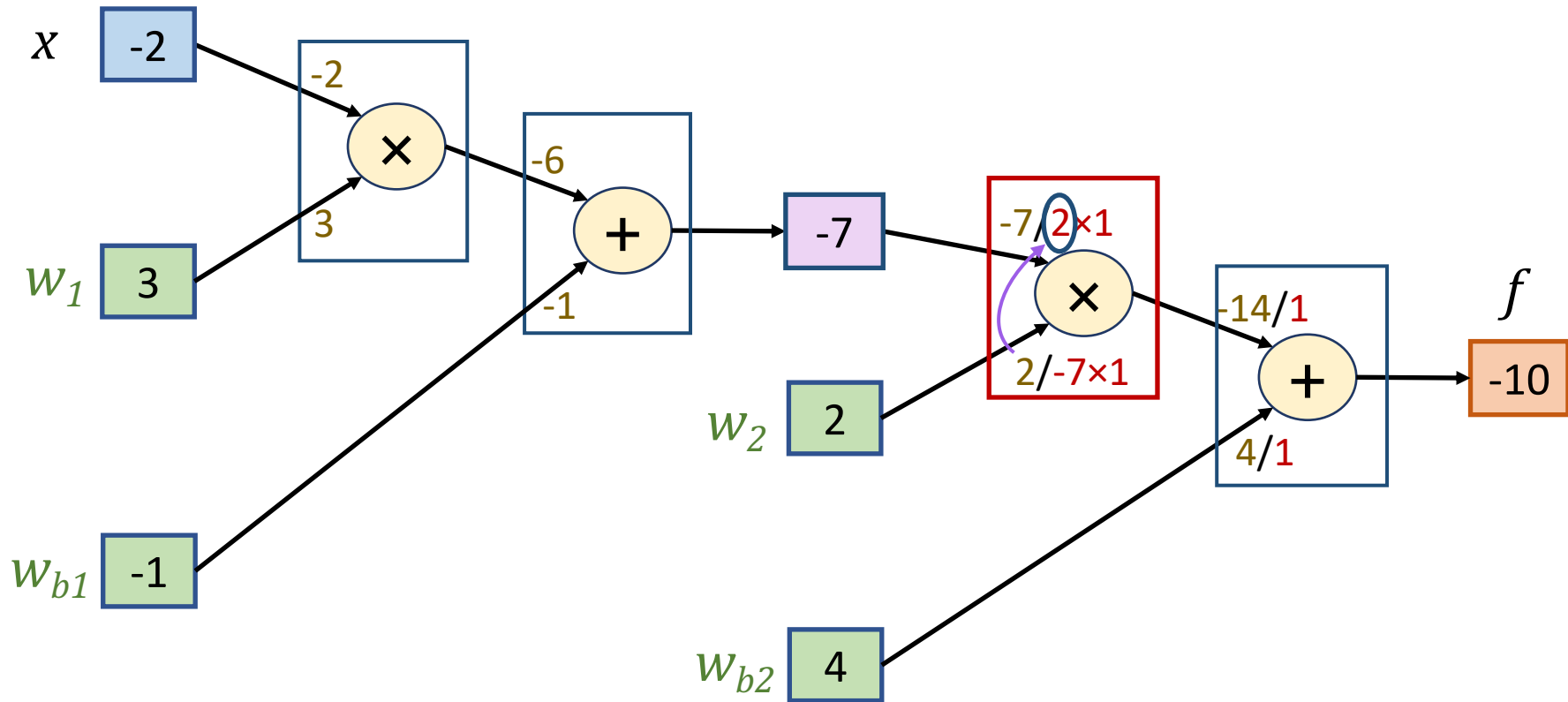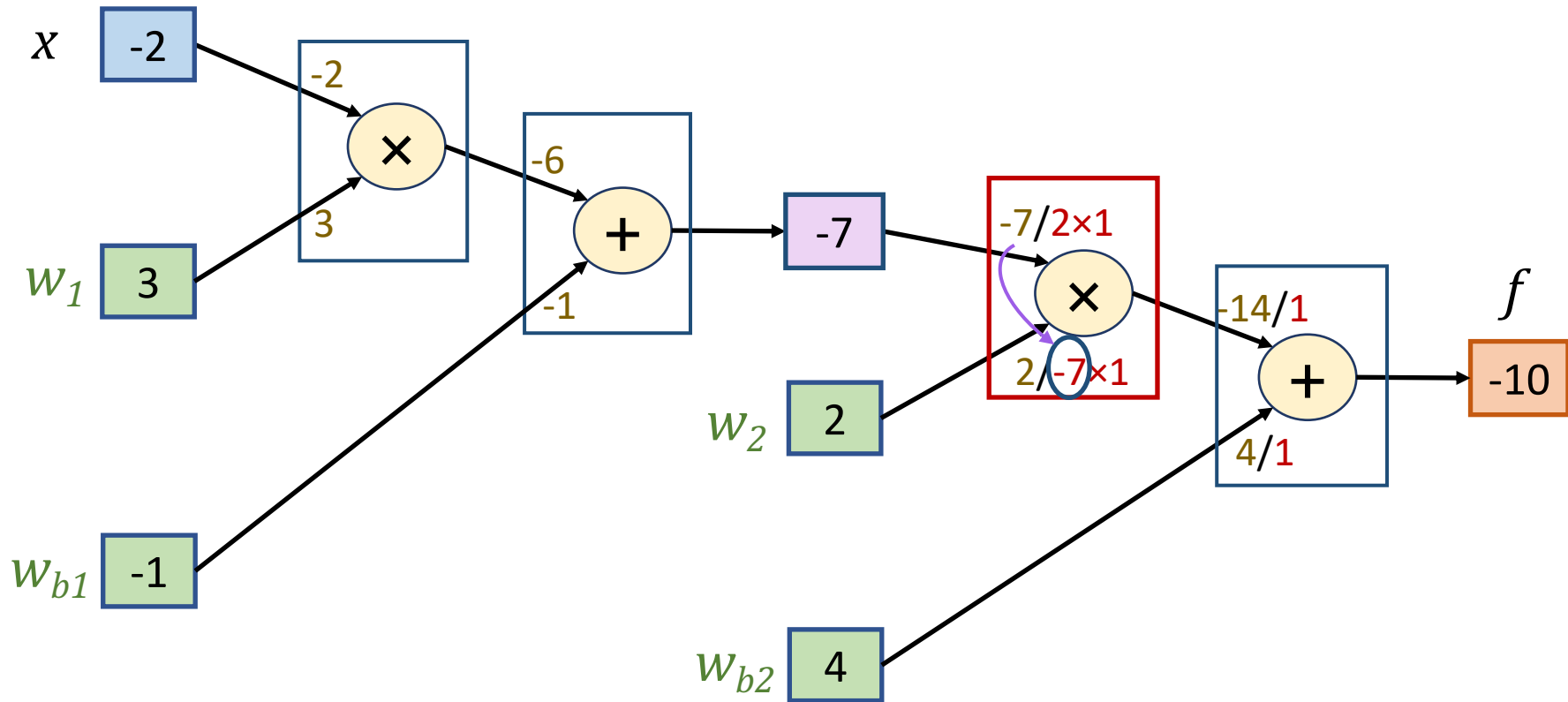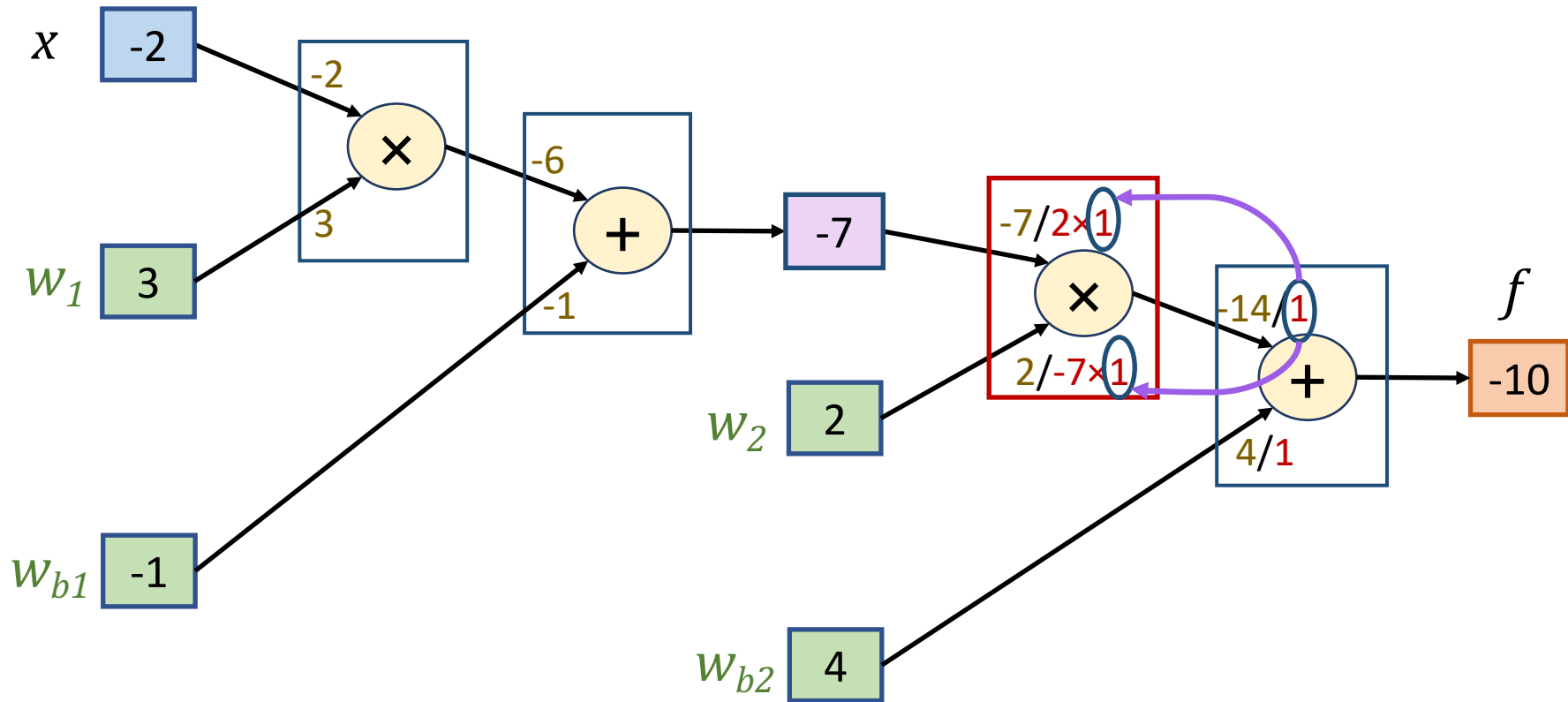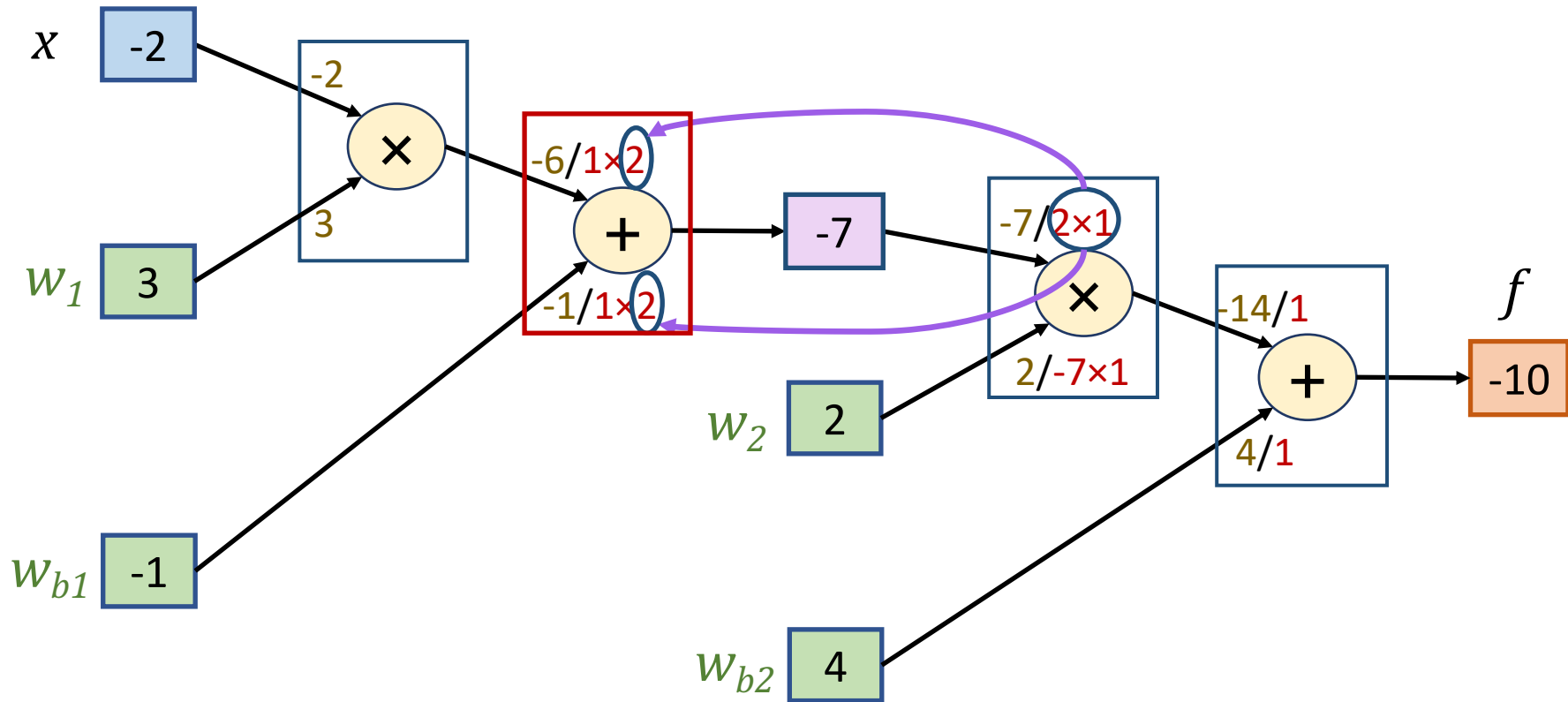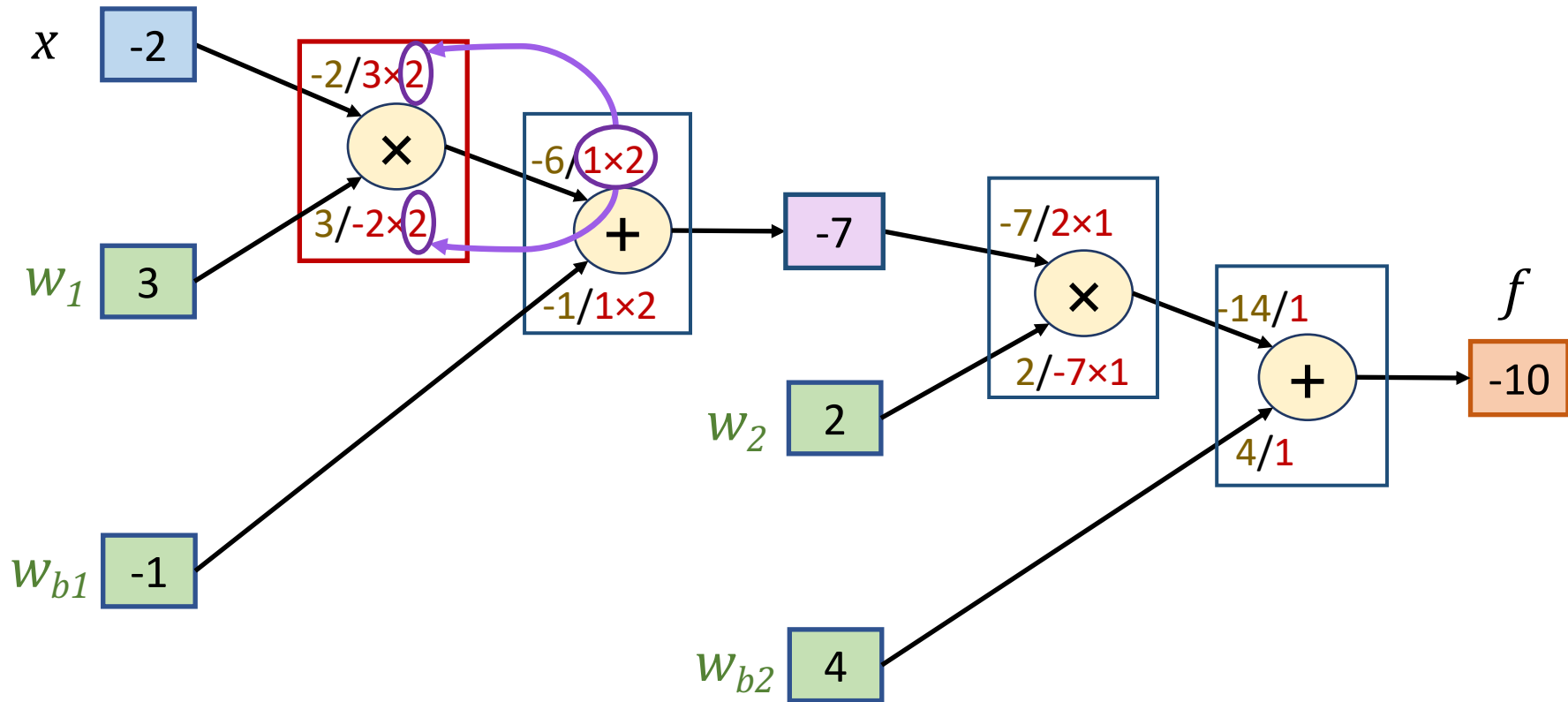
# Network Training

# Network Training

# Network Training

# Network Training

# Network Training

# Network Training

# Network Training

# Network Training

# Network Training

# Network Training

# Network Training

We can add more layers as well as our loss function

All functions must be differentiable

# Summary and Conclusion



**Forward propagation**

Network Output (*f*)

*Softmax* (*s*)

Ground truth (*g*)

| | | |
|---|---|---|
| -0.2 | *0.10* | 0 dog |
| 0.3 | *0.18* | 0 tiger |
| 0.8 | *0.31* | 0 ostrich |
| ⋮ | ⋮ | ⋮ |
| -0.1 | *0.11* | 1 Cat |
| ⋮ | ⋮ | ⋮ |
| 0.4 | *0.20* | 0 horse |

Optimisation during training

1. It makes a prediction using its current parameter set (weights).
2. We give it a feedback of how good its prediction was.
3. It updates (refines) its parameter set according to our feedback

# Summary and Conclusion

**Forward propagation**

Network Output ($f$)

*Softmax* ($s$)

Ground truth ($g$)

Deep Neural Network

| $f$ |
|------|
| -0.2 |
| 0.3 |
| 0.8 |
| ⋮ |
| -0.1 |
| ⋮ |
| 0.4 |

| $s$ |
|------|
| 0.10 |
| 0.18 |
| 0.31 |
| ⋮ |
| 0.11 |
| ⋮ |
| 0.20 |

| $g$ | |
|------|------|
| 0 | dog |
| 0 | tiger |
| 0 | ostrich |
| ⋮ | |
| 1 | Cat |
| ⋮ | |
| 0 | horse |

**Network loss**

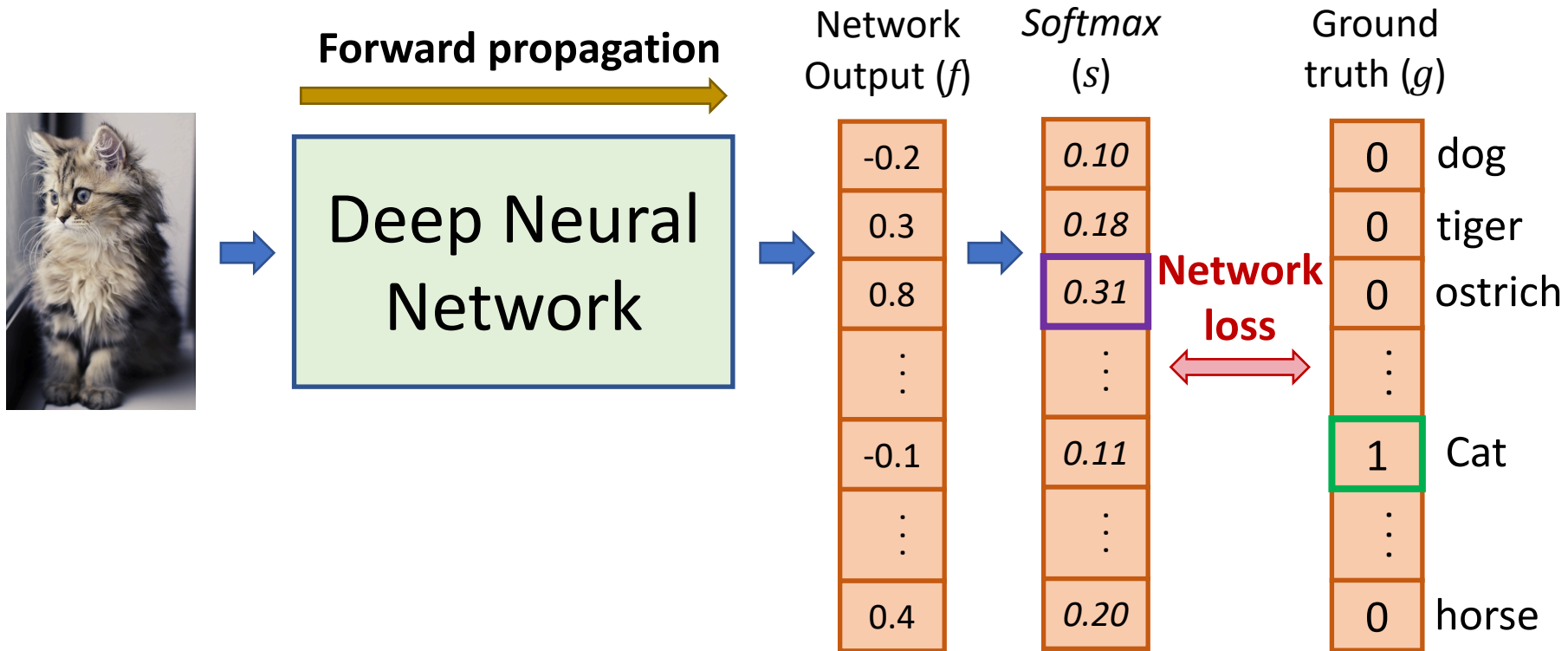## Optimisation during training

1. It makes a prediction using its current parameter set (weights).
2. We give it a feedback of how good its prediction was.
3. It updates (refines) its parameter set according to our feedback

# Summary and Conclusion



**Forward propagation**

Network Output ($f$)

*Softmax* ($s$)

Ground truth ($g$)

Deep Neural Network

| $f$ | $s$ | $g$ | |
|---|---|---|---|
| -0.2 | *0.10* | 0 | dog |
| 0.3 | *0.18* | 0 | tiger |
| 0.8 | *0.31* | 0 | ostrich |
| ⋮ | ⋮ | ⋮ | |
| -0.1 | *0.11* | 1 | Cat |
| ⋮ | ⋮ | ⋮ | |
| 0.4 | *0.20* | 0 | horse |

**Network loss**

**Backward propagation**
Using SGD and BP algorithms

## Optimisation during training

1. It makes a prediction using its current parameter set (weights).
2. We give it a feedback of how good its prediction was.
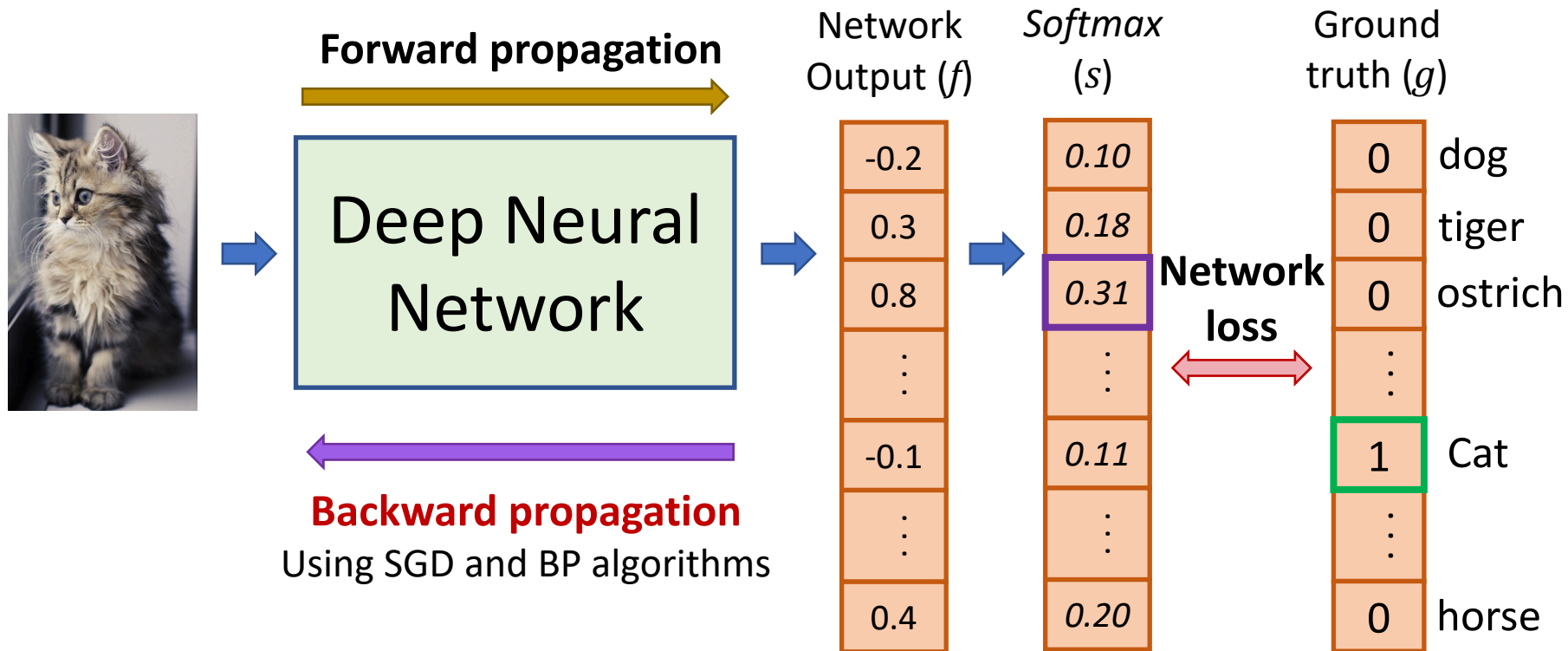3. It updates (refines) its parameter set according to our feedback

# Summary and Conclusion

- Gradient Descent is an optimization algorithm used for minimizing the network loss

- It updates the parameter set by moving in the opposite direction of the gradient of the loss function

- Back Propagation is used along with Gradient Descent to update the weights in different layers of a network

- It uses chain rule recursively to back-propagate the loss through the network

# Next part

We will discuss
different types of DNNs
and their architectures

# Thanks for your attention