

# Neural Networks

# Neural networks

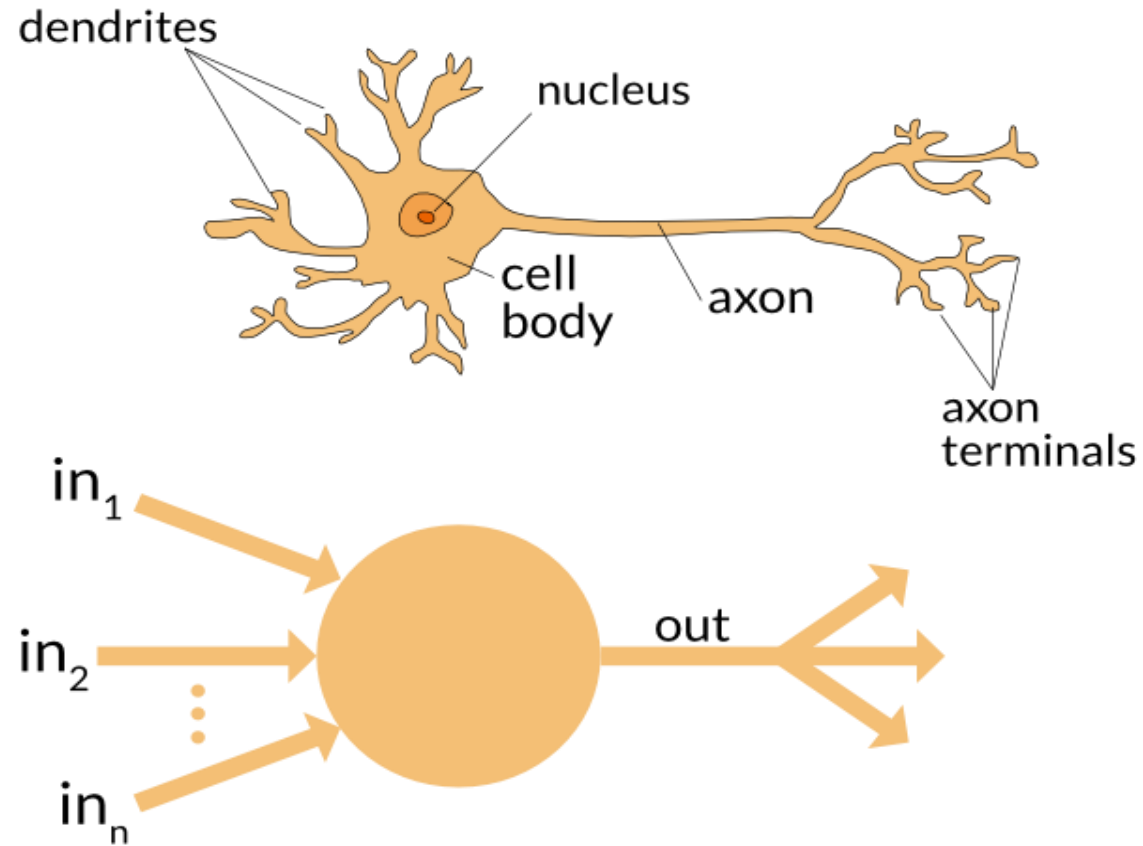
**Goal:** Build an intuition into how neural networks works and get familiar with the associated vocabulary.

## Program:

- Intuition to neural networks
- The anatomy of a network
- The training process (gradient descent)
- Backpropagation

# Intuition

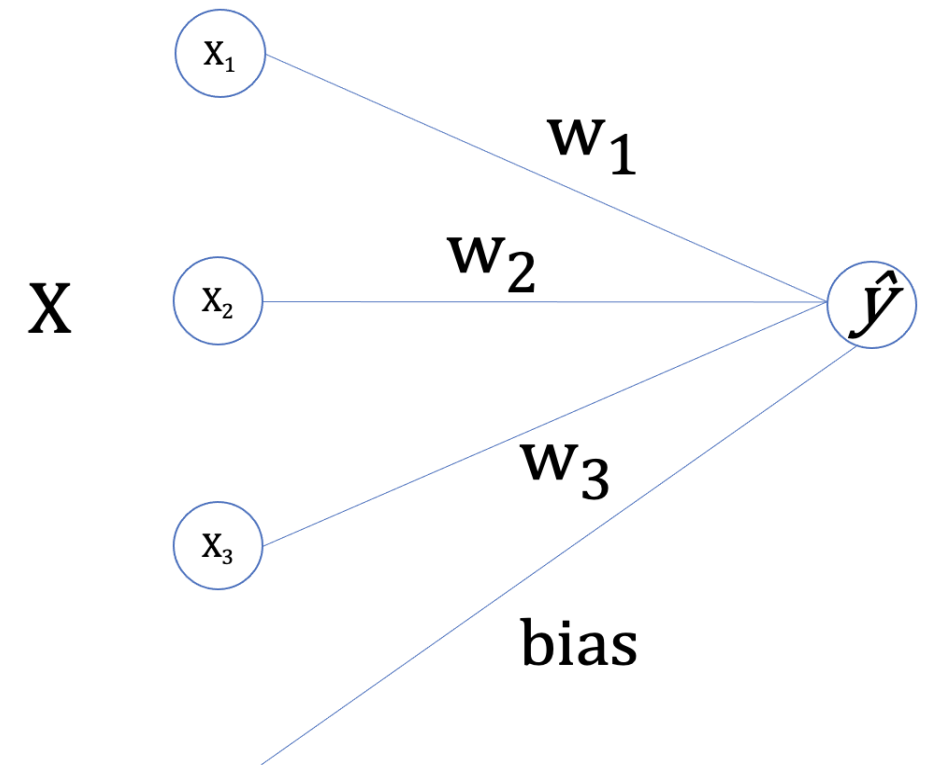
# Neural Networks Intuition



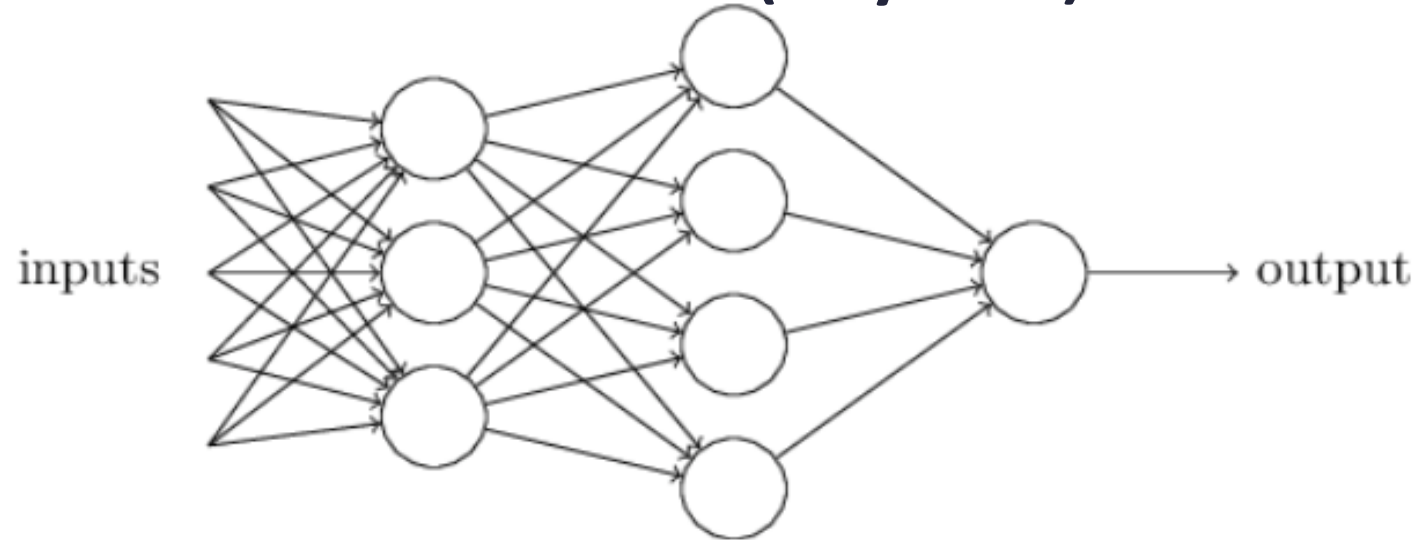
# Anatomy of the network

# Anatomy of the network (neurons)

Inputs (vector  $\mathbf{x}$ ) are combined  
with weights (vector  $\mathbf{w}$ ),  
a bias (value  $b$ )  
and a non-linear activation function ( $\sigma$ )  
to produce an output value,  
i.e.  $\text{output} = \sigma(\mathbf{w}^T \mathbf{x} + b)$



# Anatomy of the network (layers)



**Layer:** column of neurons stacked together that can receive the same inputs.

**Hidden layer:** intermediate layers between inputs and outputs.

**Deep neural network:** a neural network that contains many hidden layers, and can therefore provide solutions to more complicated and subtle decision problems.

# The training process



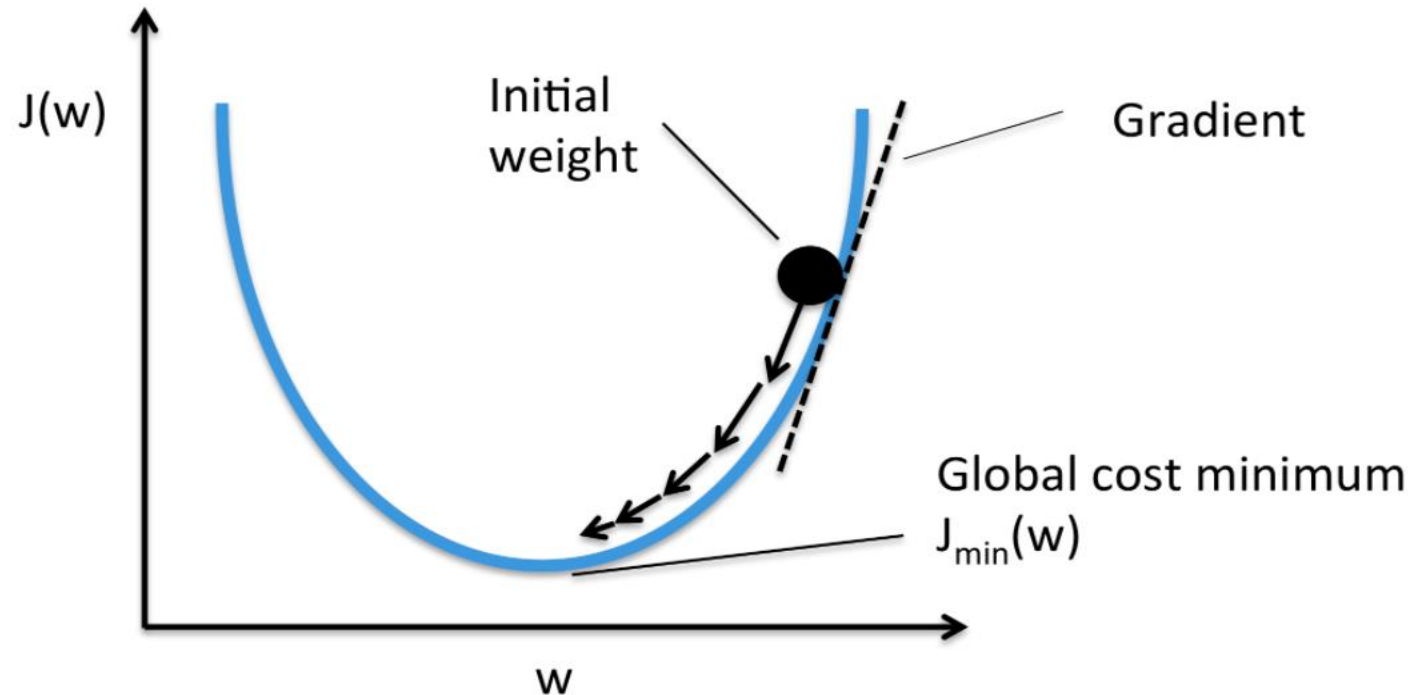
# Training

- Training is the process of tuning the weights  $\mathbf{w}$  in a network.
- Usually, to train a network we require input data  $\mathbf{X}$  and corresponding target labels  $\mathbf{y}$ .
- We attempt to use the network  $f$  to predict make predictions, i.e.  $\hat{\mathbf{y}} = f(\mathbf{X}, \mathbf{w})$ .
- The weights should minimise a **cost or loss function**  $J$ , e.g.
  - $J = \mathbf{y} - \hat{\mathbf{y}}$
  - $J = \frac{1}{2} (\mathbf{y} - \hat{\mathbf{y}})^2$

# Gradient descent

$$w' = w - \eta \frac{\partial J(w)}{\partial w}$$

A process of following the gradients of the error function towards a minimum value.



# Backpropagation

# Backpropagation

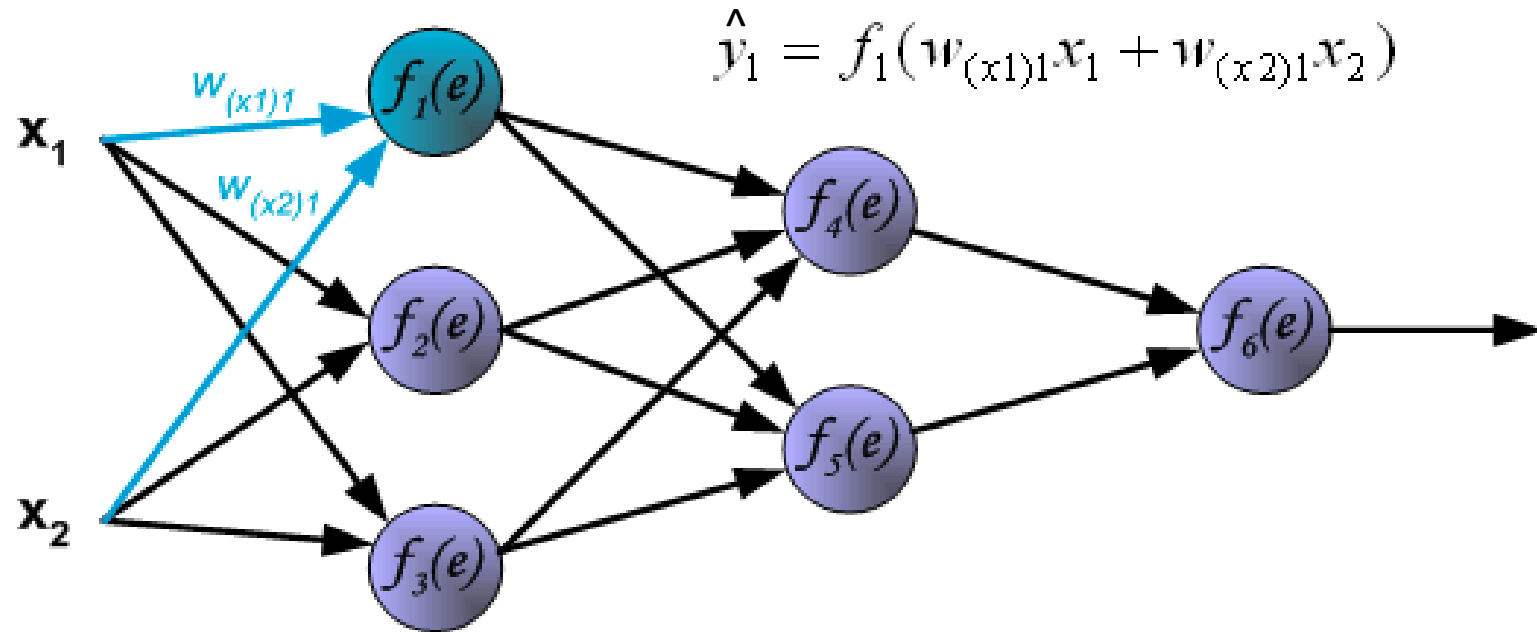
We'll now examine backpropagation, a fast algorithm for computing the weight gradients which is based on the chain-rule.



# Weight updates

- We want to find the best weights!
- Update weights using gradient descent

- $w'_{(x1)1} = w_{(x1)1} - \eta \frac{\partial loss}{\partial w_{(x1)1}}$

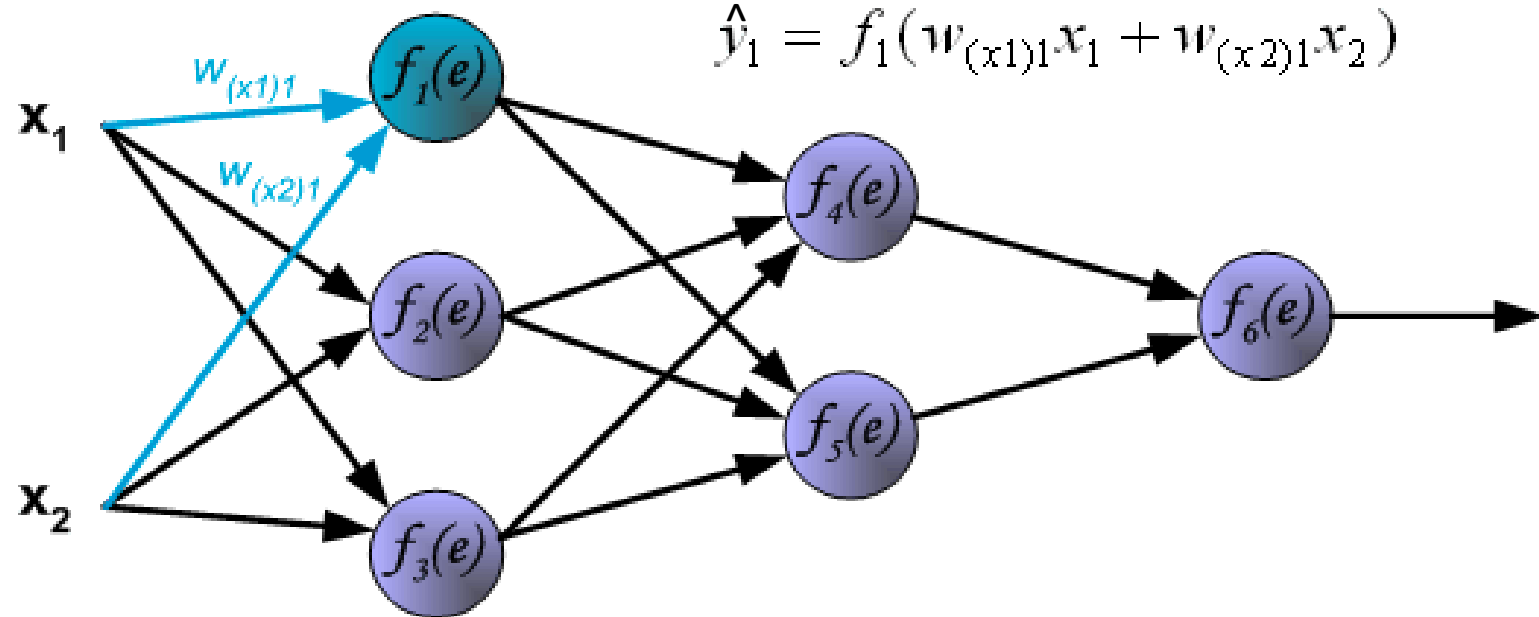


# Weight updates

- We find  $\frac{\partial \text{loss}}{\partial w_{(x1)1}}$  using the chain rule

**We need to compute these gradients**

- $\frac{\partial \text{loss}}{\partial w_{(x1)1}} = \frac{\partial \text{loss}}{\partial f_1(e)} \frac{\partial f_1(e)}{\partial e} \frac{\partial e}{\partial w_{(x1)1}}$ , where  $e = w_{(x1)1}x_1 + w_{(x2)1}x_2$



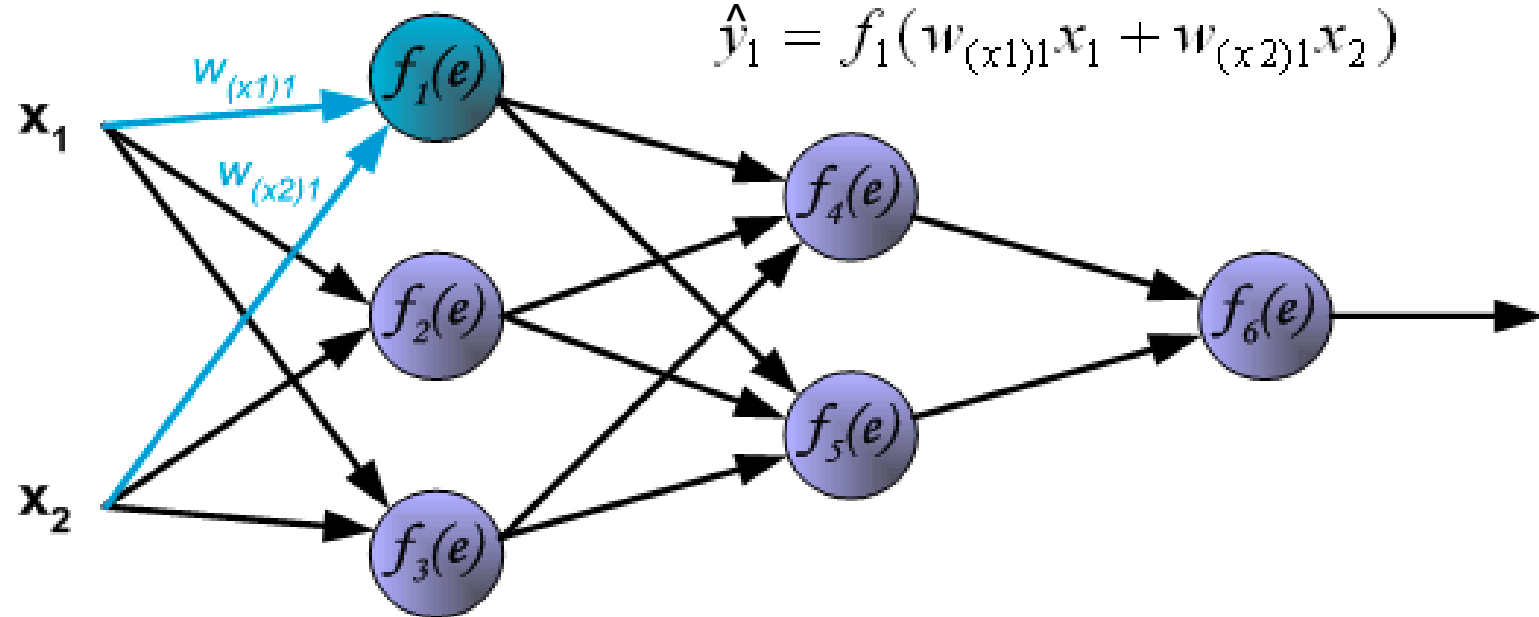
# Weight updates

- We find  $\frac{\partial \text{loss}}{\partial w_{(x1)1}}$  using the chain rule

**We need to compute these gradients**

- $\frac{\partial \text{loss}}{\partial w_{(x1)1}} = \frac{\partial \text{loss}}{\partial f_1(e)} \frac{\partial f_1(e)}{\partial e} \frac{\partial e}{\partial w_{(x1)1}}$ , where  $e = w_{(x1)1}x_1 + w_{(x2)1}x_2$

$$\frac{\partial e}{\partial w_{(x1)1}} = x_1$$



# Weight updates

- We find  $\frac{\partial loss}{\partial w_{(x1)1}}$  using the chain rule

**We need to compute these gradients**

- $\frac{\partial loss}{\partial w_{(x1)1}} = \frac{\partial loss}{\partial f_1(e)} \frac{\partial f_1(e)}{\partial e} x_1$

- This is the derivative of our activation function, typically we choose one that is easy to compute, e.g. the sigmoid

$$\sigma'(e) = \sigma(e)(1 - \sigma(e))$$



# Weight updates

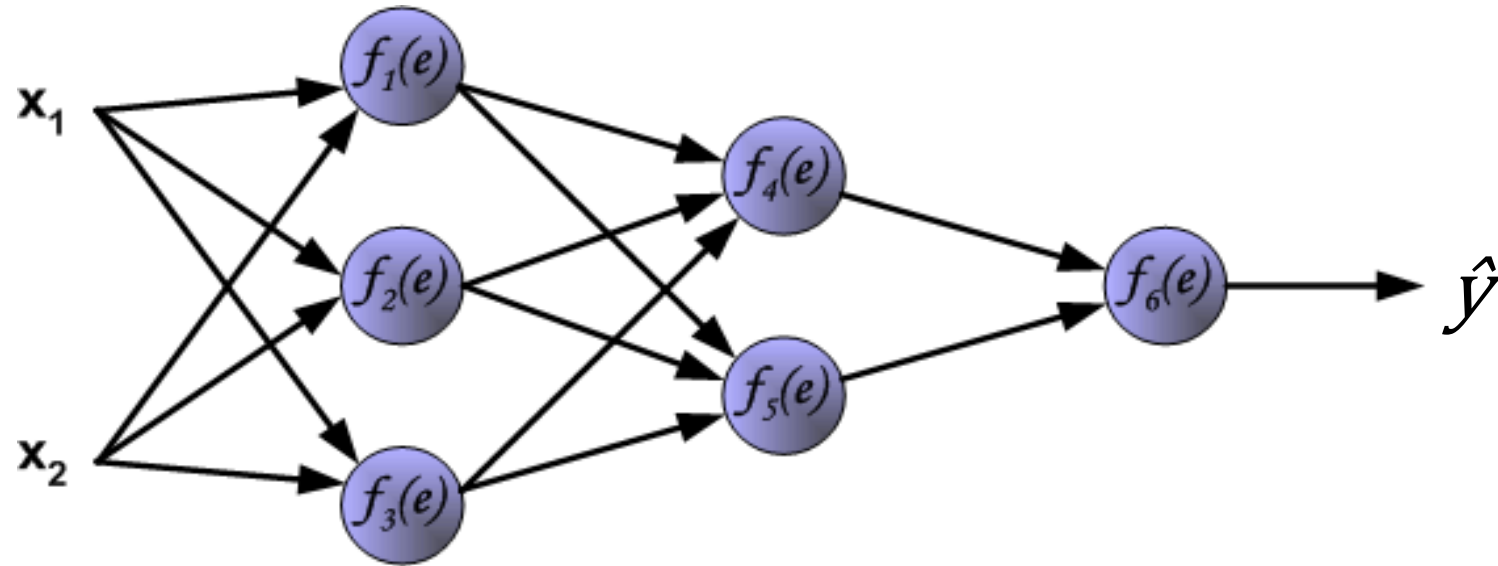
- We find  $\frac{\partial loss}{\partial w_{(x1)1}}$  using the chain rule

- $\frac{\partial loss}{\partial w_{(x1)1}} = \frac{\partial loss}{\partial f_1(e)} \frac{\partial f_1(e)}{\partial e} x_1$

- We compute this derivative by backpropagating the error through the network.

**We need to compute these gradients**

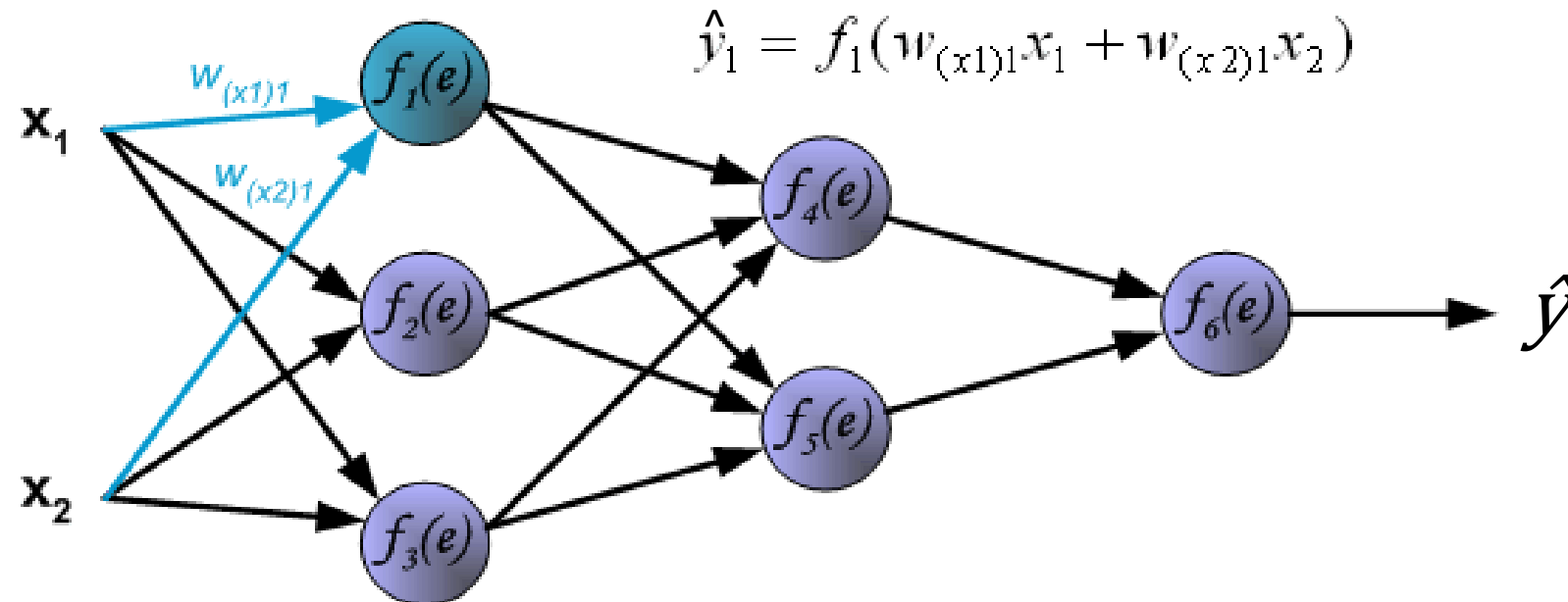
# Epoch



- **Forward pass:** a data sample is passed forward through the network to determine a prediction.
- **Backward pass:** recursively compute the error backwards from the last layer following the chain-rule and update the weights w.r.t. the known target output.
- **Epoch:** training the neural network with all the training data for one cycle.

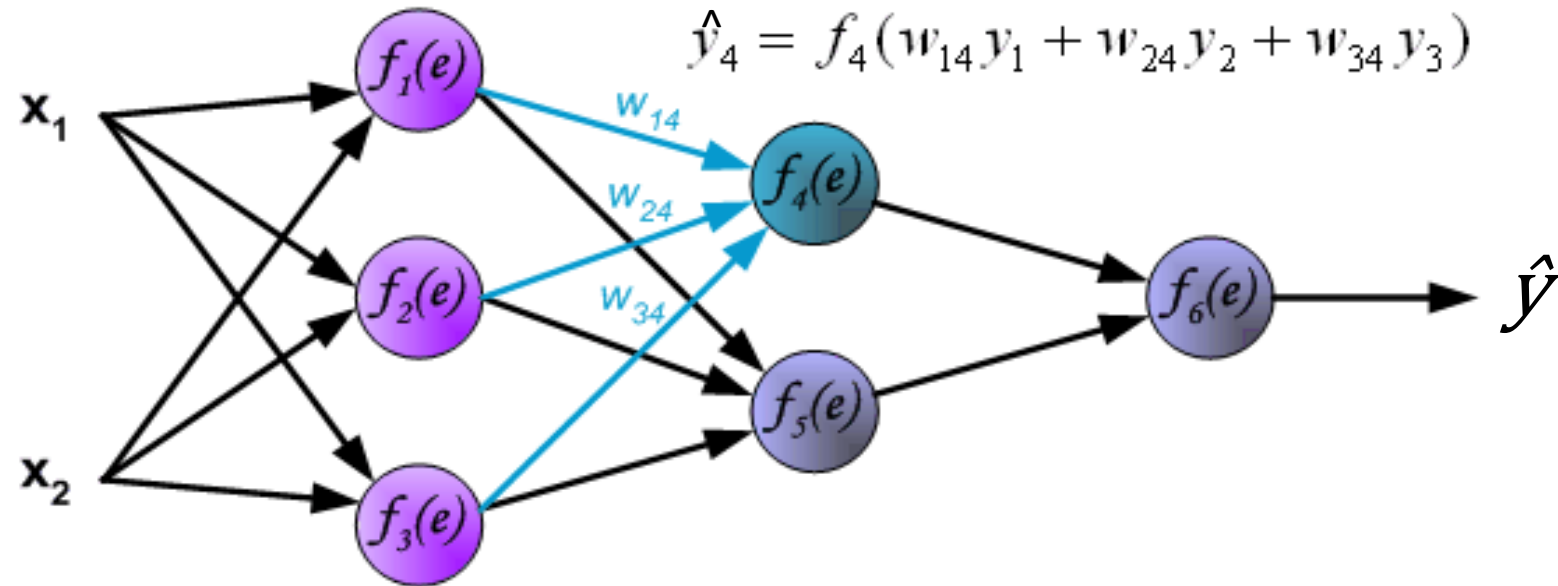
# Forward pass

Feed data through the network.



# Forward pass

Feed data through the network.



# Compute the error

For example,

$$\text{loss} = \frac{1}{2} (\mathbf{y} - \hat{\mathbf{y}})^2$$

Then we can compute the error rate  $\delta = \frac{\partial \text{loss}}{\partial \hat{\mathbf{y}}}$

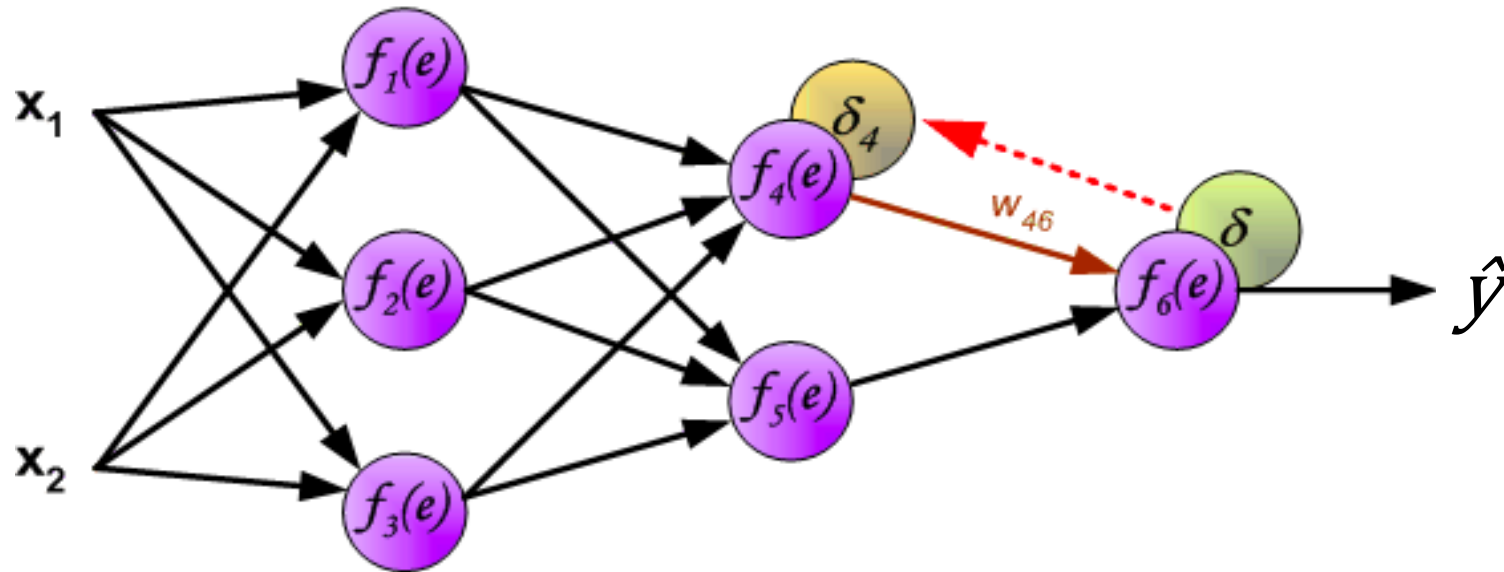
$$\delta = \frac{\partial \text{loss}}{\partial \hat{\mathbf{y}}} = -(\mathbf{y} - \hat{\mathbf{y}})$$

# Backpropagation

## Local error contribution

$$\delta_4 = \frac{\partial \text{loss}}{\partial f_4} = \frac{\partial \text{loss}}{\partial f_6} \frac{\partial f_6(e)}{\partial e} \frac{\partial e}{\partial f_4} = \delta \frac{\partial f_6(e)}{\partial e} w_{46}$$

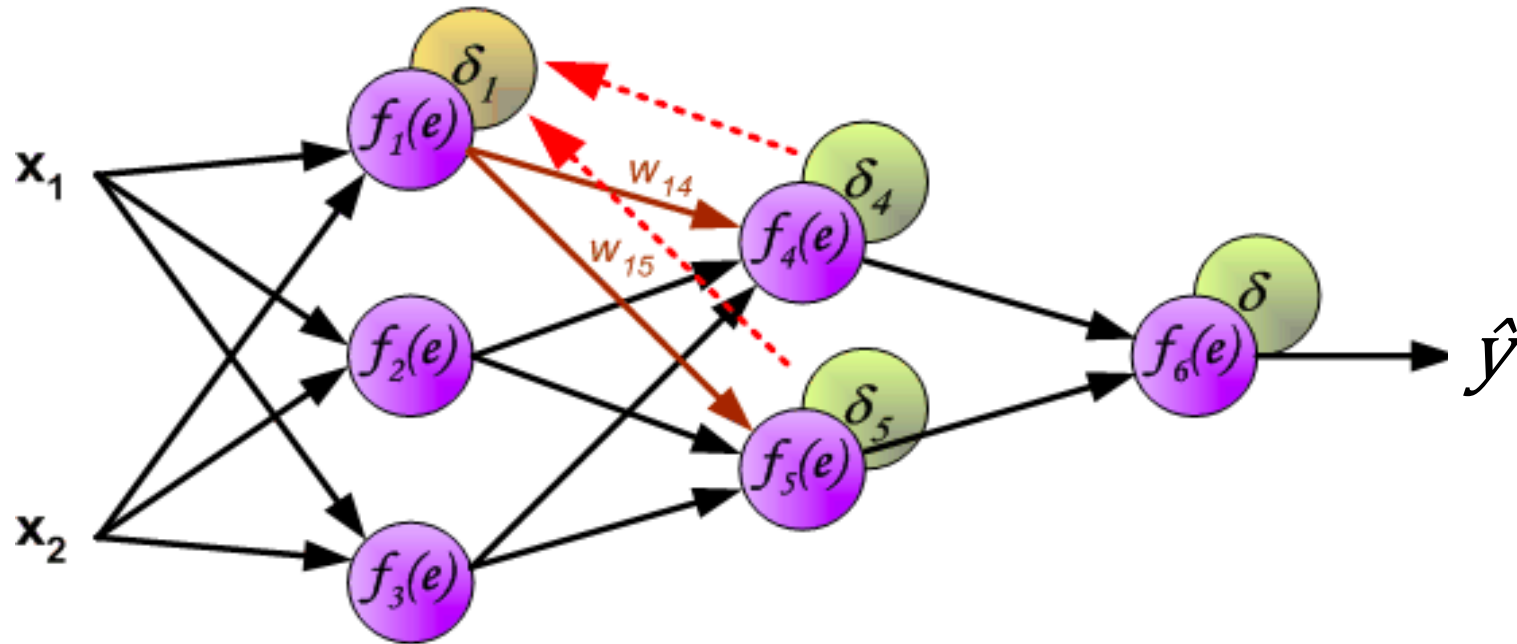
where  $\hat{y} = f_6$  &  $e = w_{46}f_4 + w_{56}f_6$



# Backpropagation

Local error contribution

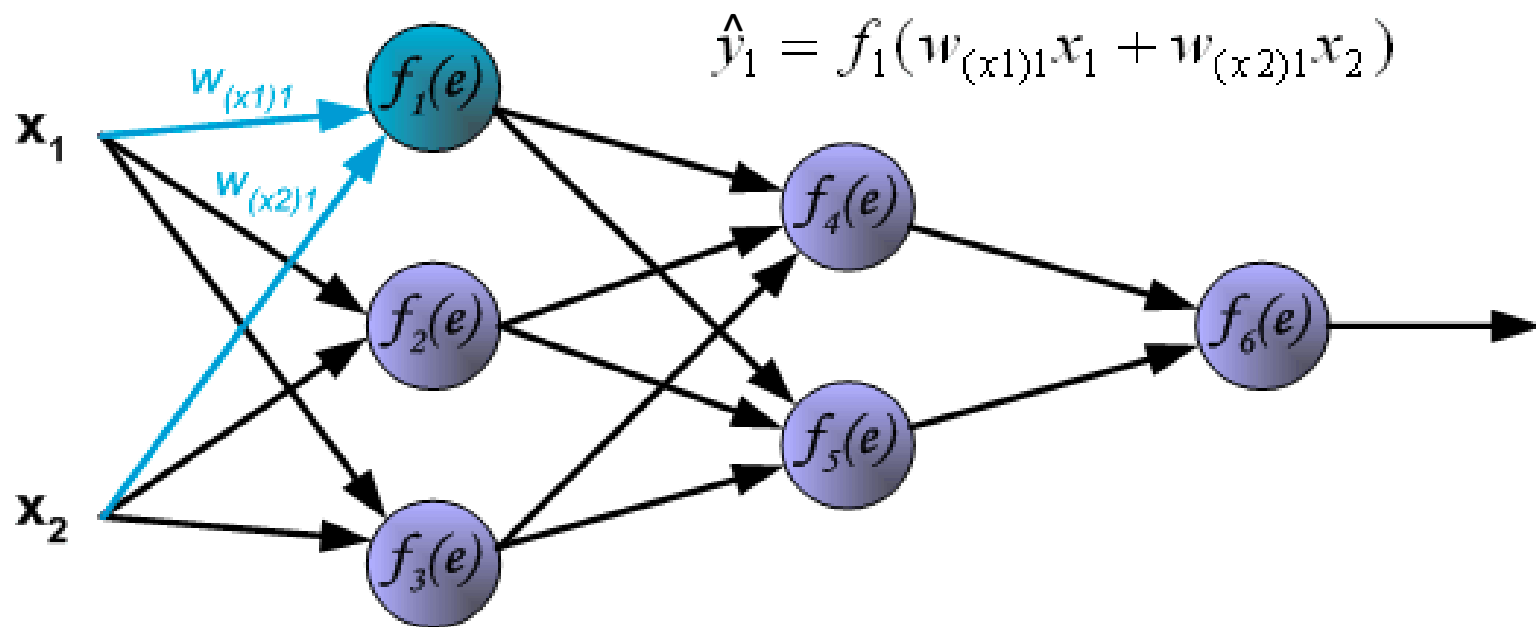
$$\delta_1 = \frac{\partial \text{loss}}{\partial f_1} = \frac{\partial \text{loss}}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial e} \frac{\partial e}{\partial f_4} = \delta_4 \frac{\partial f_4(e)}{\partial e} w_{14} + \delta_5 \frac{\partial f_5(e)}{\partial e} w_{15}$$



# Backpropagation

Previously we saw that

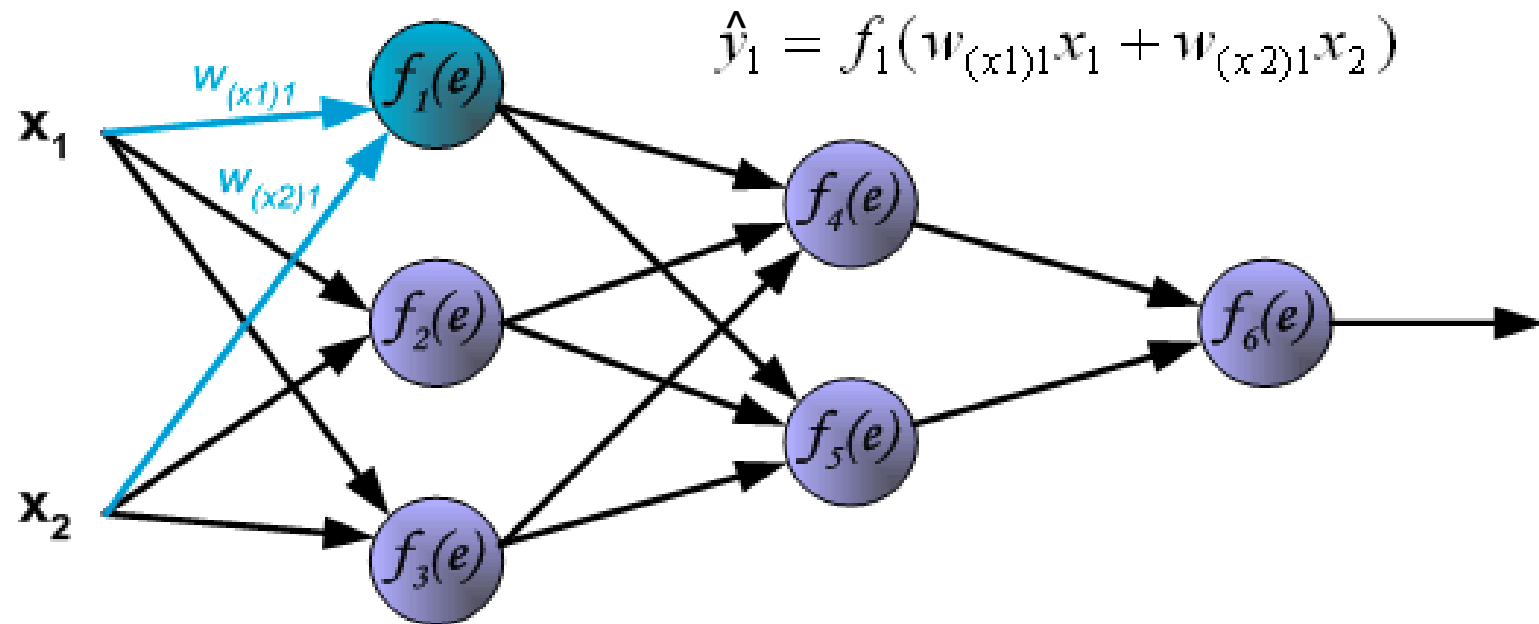
$$\frac{\partial \text{loss}}{\partial w_{(x1)1}} = \frac{\partial \text{loss}}{\partial f_1(e)} \frac{\partial f_1(e)}{\partial e} w_{(x1)1}$$





# Backpropagation

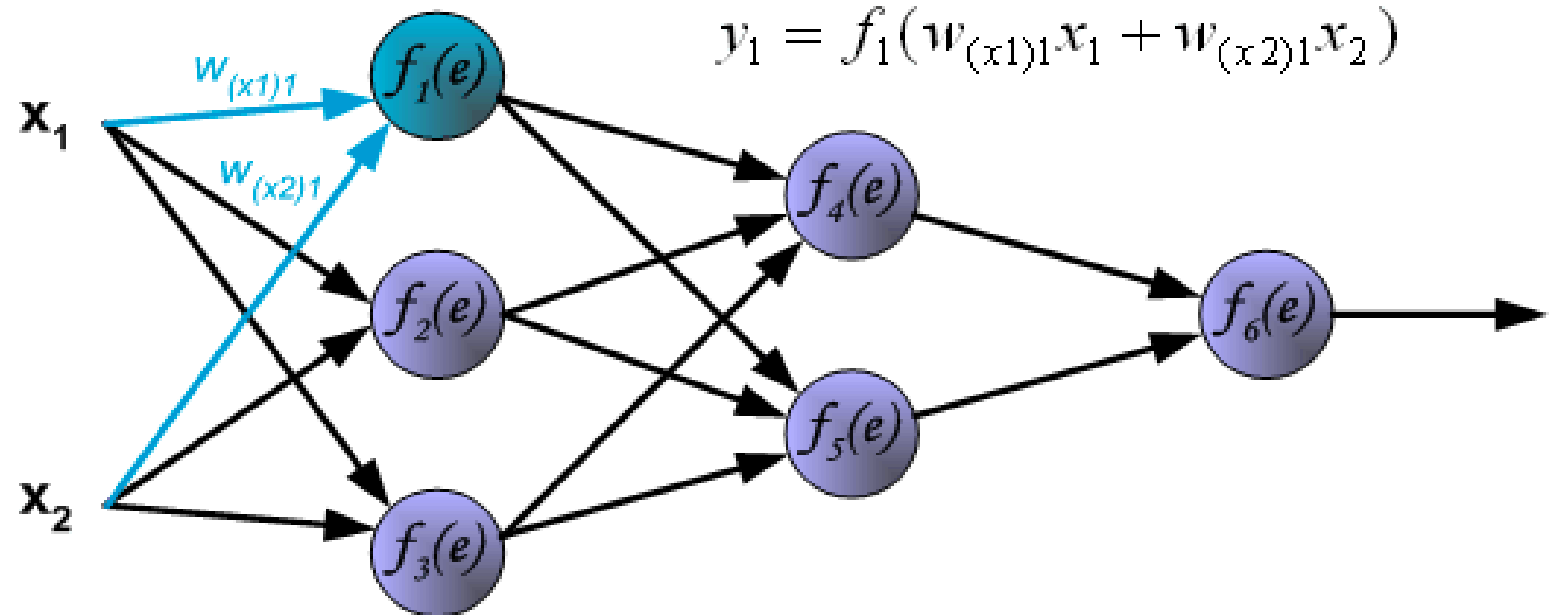
$$\text{Thus, } \frac{\partial \text{loss}}{\partial w_{(x1)1}} = \delta_1 \frac{\partial f_1(e)}{\partial e} x_1$$



# Backpropagation

We can then optimise weight using gradient descent, e.g.

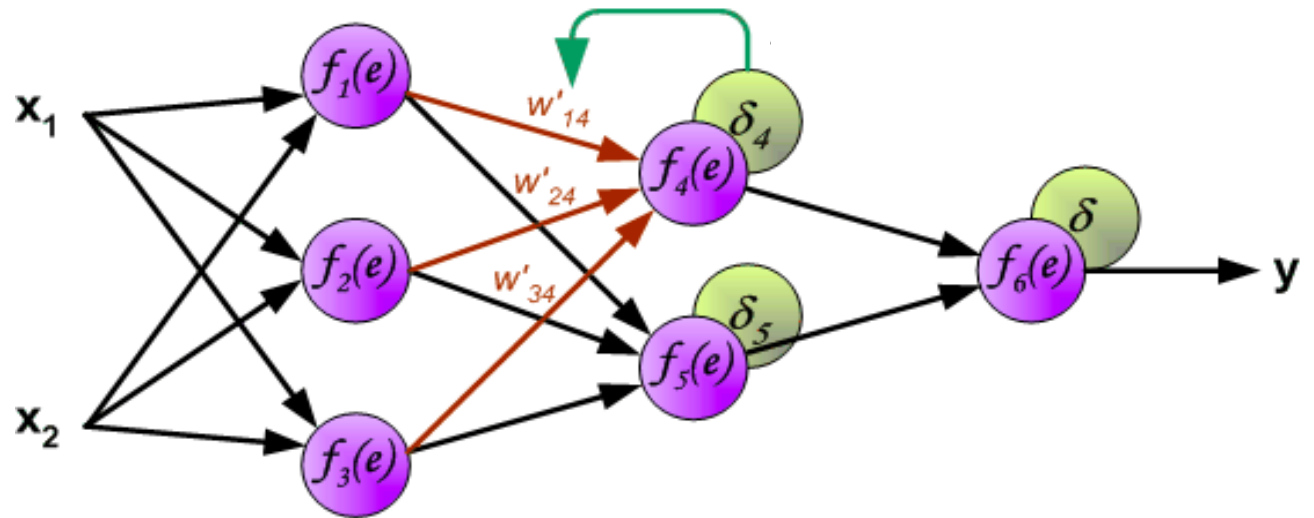
$$w'_{(x1)1} = w_{(x1)1} - \eta \delta_1 \frac{\partial f_1(e)}{\partial e} x_1$$



# Backpropagation

We can then optimise weight using gradient descent, e.g.

$$w'_{14} = w_{14} - \eta \delta_4 \frac{\partial f_4(e)}{\partial e} f_1$$



# Summary

# Summary

In this section we have covered

- Intuition to neural networks
- The anatomy of a network
- The training process (gradient descent)
- Backpropagation

In the next exercise we will build a MLP to solve the XOR problem

# XOR perceptron

# XOR

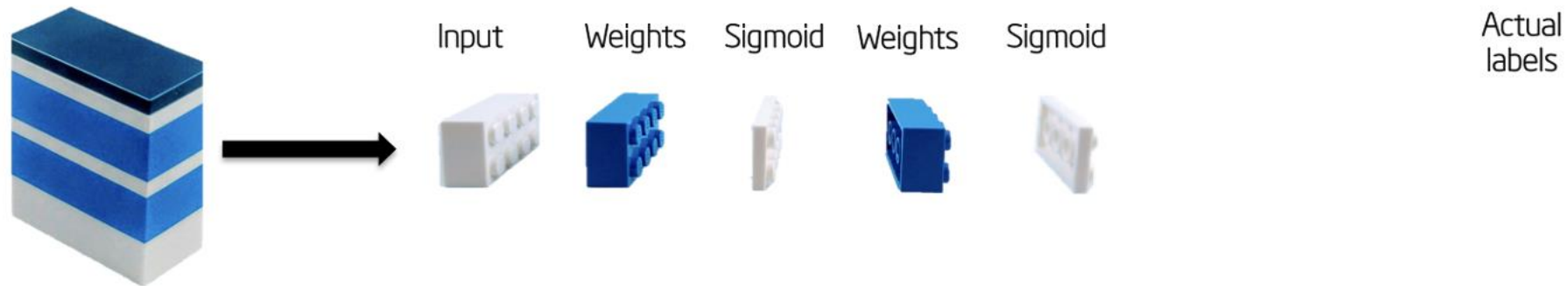
A XOR gate (or exclusive OR) is a logic gate that returns true (or 1) when one, and only one, of the inputs of the gate is true. Otherwise it returns false.

A	B	A XOR B
0	0	0
0	1	1
1	0	1
1	1	0

Whereas other logical operators (AND, OR) can be modeled with a single layer perceptron, the XOR operator is more complex and must be modeled with a multi-layer perceptron.

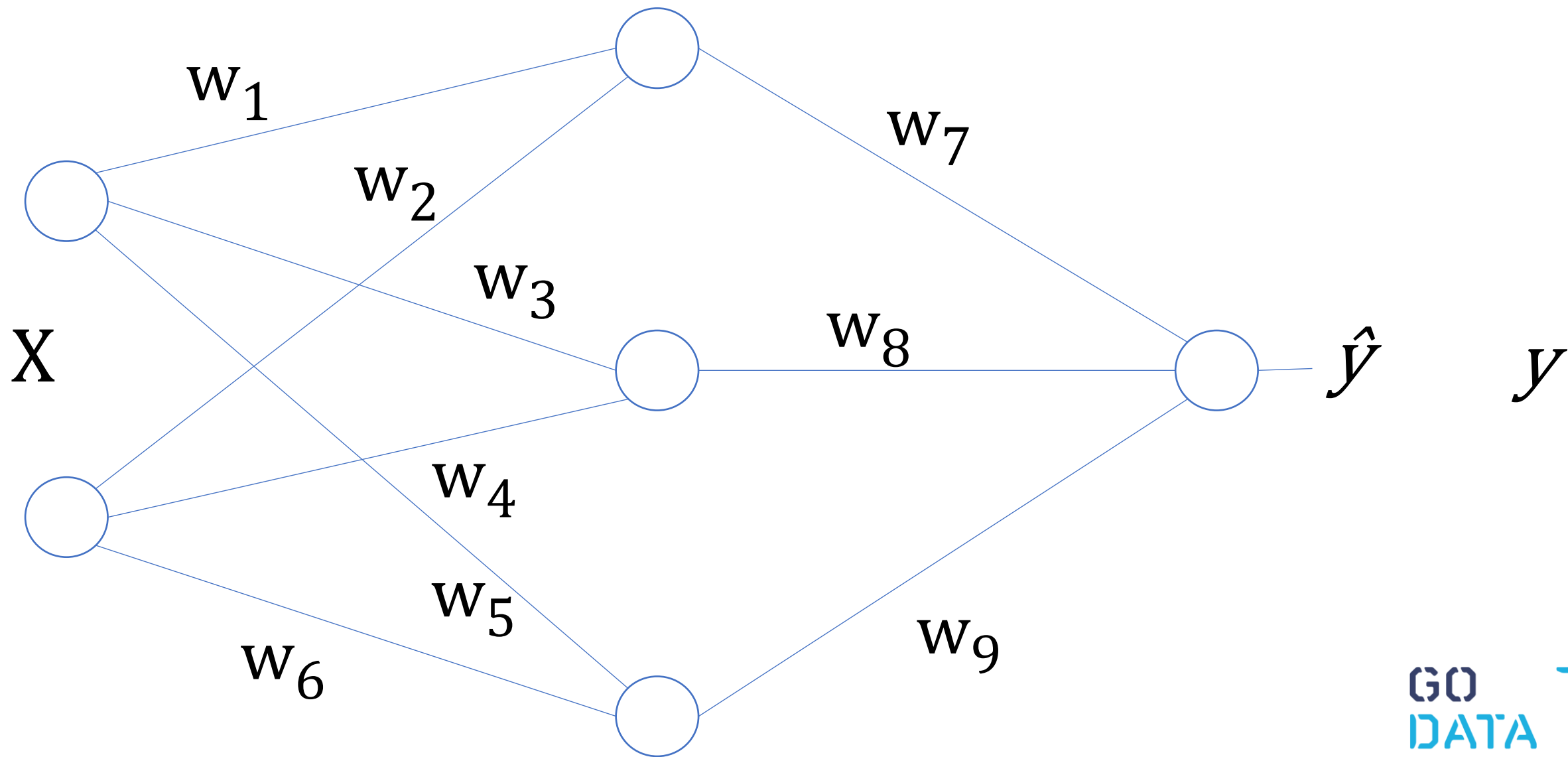
# MLP

The MLP applies layers of transformations to our input. We compare the output of the model to our actual labels.

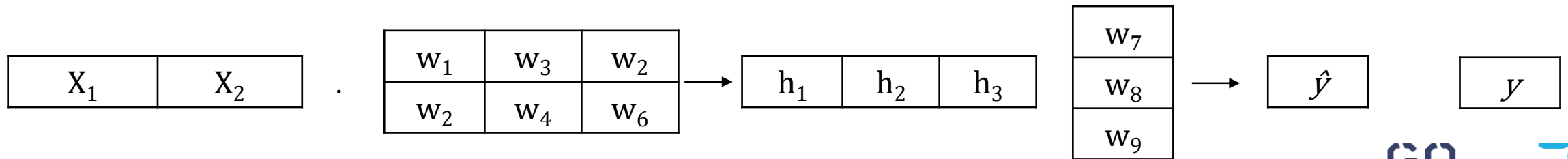
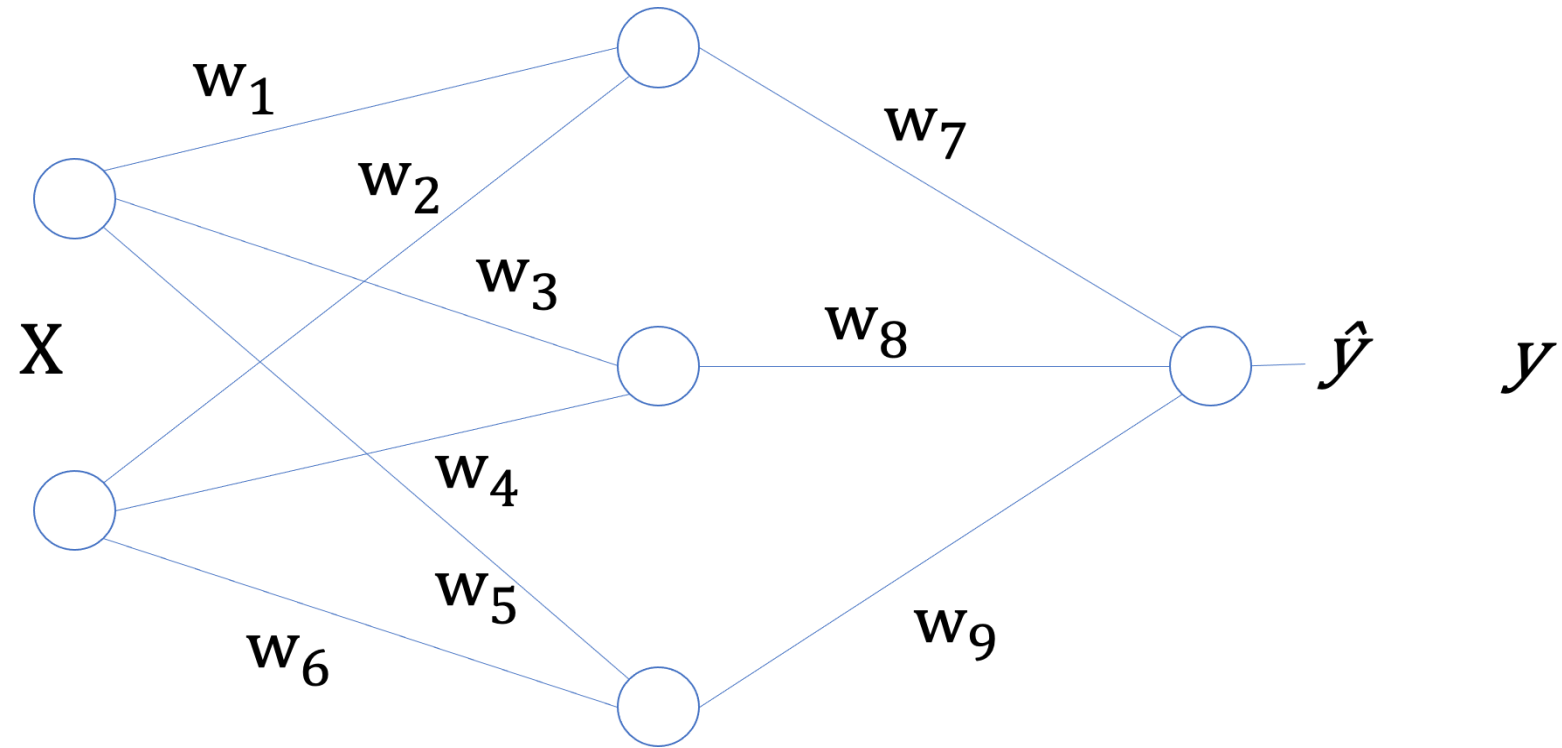


The next slide shows a schematic diagram of the MLP we will be creating





# XOR MLP



# Questions?

GO   
DATA  
DRIVEN