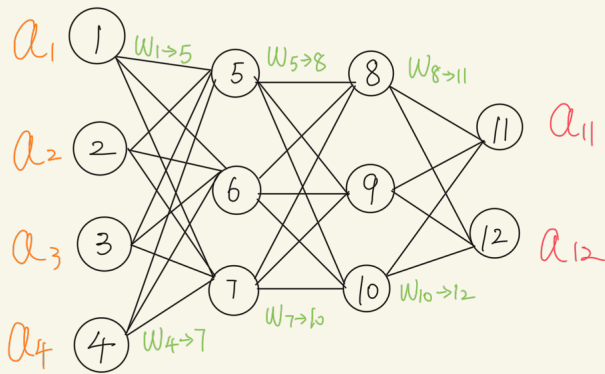


# Feed Forward Neural Networks



input dimension: 4

output dimension: 2

hidden layers: [3, 3]

Computation feeds forward:

- set inputs
- loop through layers in the network
  - compute each node's activation (and store it!)
  - uses prev. layer's activations

Why we use the hidden layers?

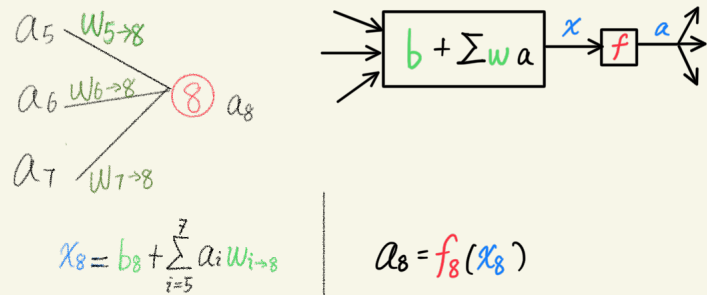
- linear activations: theorem 1

Does Nothing (nothing difference with one single neuron)

If we compute some linear functions with a neuron, then pass it through some other linear functions, we could have achieved same thing with one single neuron (with diff. weights).

- non-linear activations: theorem 2

Computes Anything

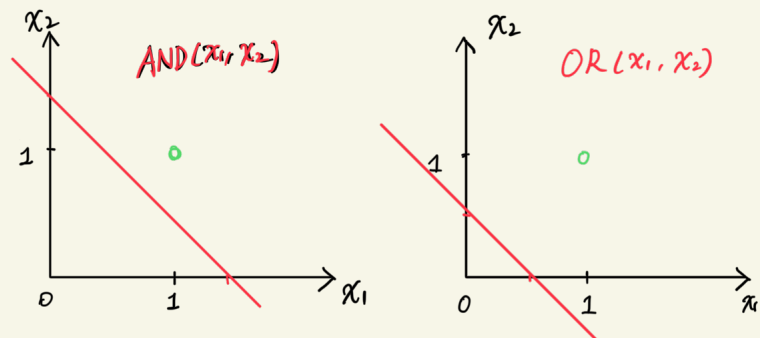


## Activation Functions

linear: $a = x$		$\frac{da}{dx} = 1$
sigmoid: $a = \frac{1}{1 + e^{-x}}$		$\frac{da}{dx} = a(1-a)$
tanh: $a = \frac{1 + e^{-2x}}{1 - e^{-2x}}$		$\frac{da}{dx} = 1 - a^2$
ReLU: $a = \text{Max}(0, x)$		$\frac{da}{dx} = \begin{cases} 0 & a=0 \\ 1 & a>0 \end{cases}$

Theorem 1: Let  $f(\vec{x})$  and  $g(\vec{x})$  be linear functions, then  $f(g(\vec{x}))$  is linear too.

Theorem 2: If we can build the functions AND/OR/NOT, then we can compose them to represent any Boolean function.



$$\text{AND}(0,0) = 0$$

$$\text{AND}(1,0) = 0$$

$$\text{AND}(0,1) = 0$$

$$\text{AND}(1,1) = 1$$

$$\text{OR}(0,0) = 0$$

$$\text{OR}(1,0) = 1$$

$$\text{OR}(0,1) = 1$$

$$\text{OR}(1,1) = 1$$

# Gradient Descent Training

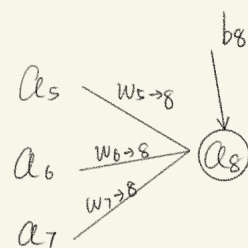
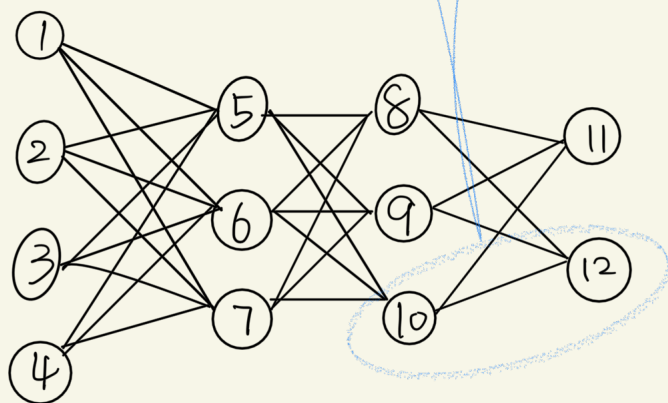
For each data point, the loss function needs to aggregate over output dimensions.

$$L^j(\theta, \text{data}^j) = \sum_{o \in \text{output}} (t_o^j - a_o^j)^2 \quad \text{Sum over outputs}$$

$$L(\theta, \text{dataset}) = \frac{1}{N} \sum_{j \in \mathcal{H}} L^j(\theta, \text{data}^j) \quad \text{Mean over data}$$

$$\theta = \begin{bmatrix} w_{1 \rightarrow 5} \\ w_{1 \rightarrow 6} \\ w_{1 \rightarrow 7} \\ \vdots \\ w_{10 \rightarrow 12} \\ b_5 \\ b_6 \\ \vdots \\ b_{11} \\ b_{12} \end{bmatrix}$$

$$\nabla L = \begin{bmatrix} \partial L / \partial w_{1 \rightarrow 5} \\ \partial L / \partial w_{1 \rightarrow 6} \\ \partial L / \partial w_{1 \rightarrow 7} \\ \vdots \\ \partial L / \partial w_{10 \rightarrow 12} \\ \partial L / \partial b_5 \\ \vdots \\ \partial L / \partial b_{11} \\ \partial L / \partial b_{12} \end{bmatrix}$$



$$x_8 = b_8 + \sum_{i=5}^7 a_i w_{i \rightarrow 8}$$