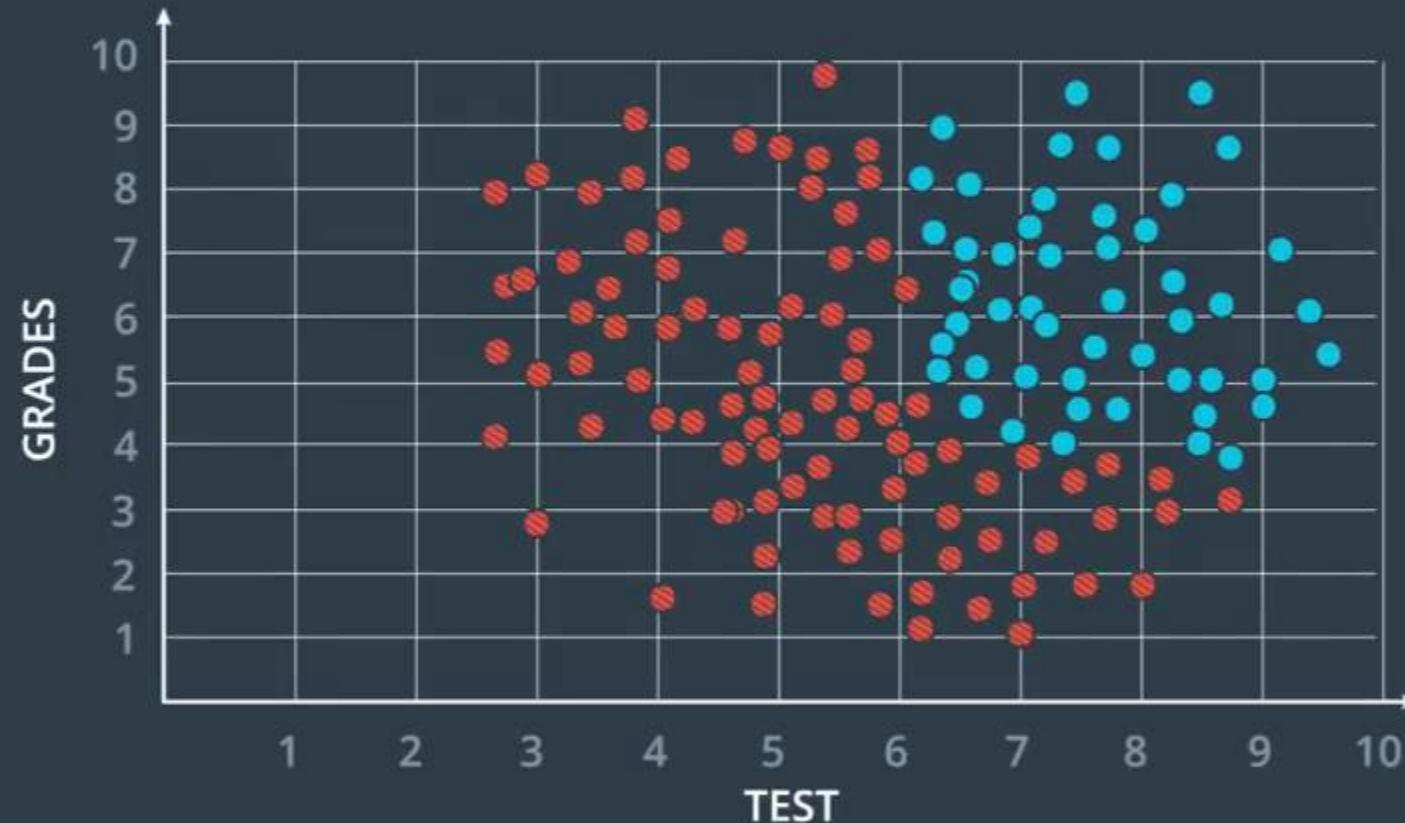


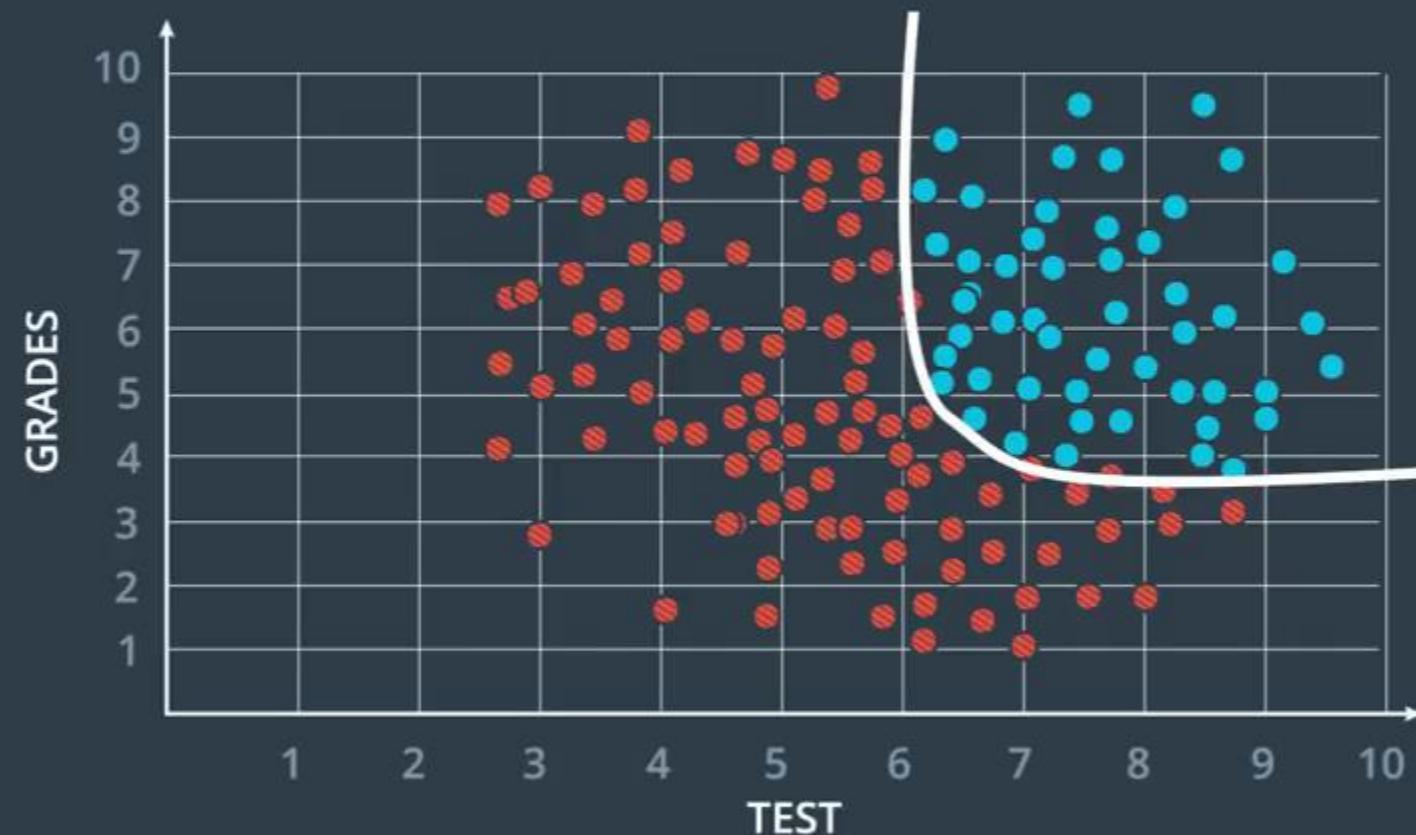
Neural Network

by Hassan Badawy

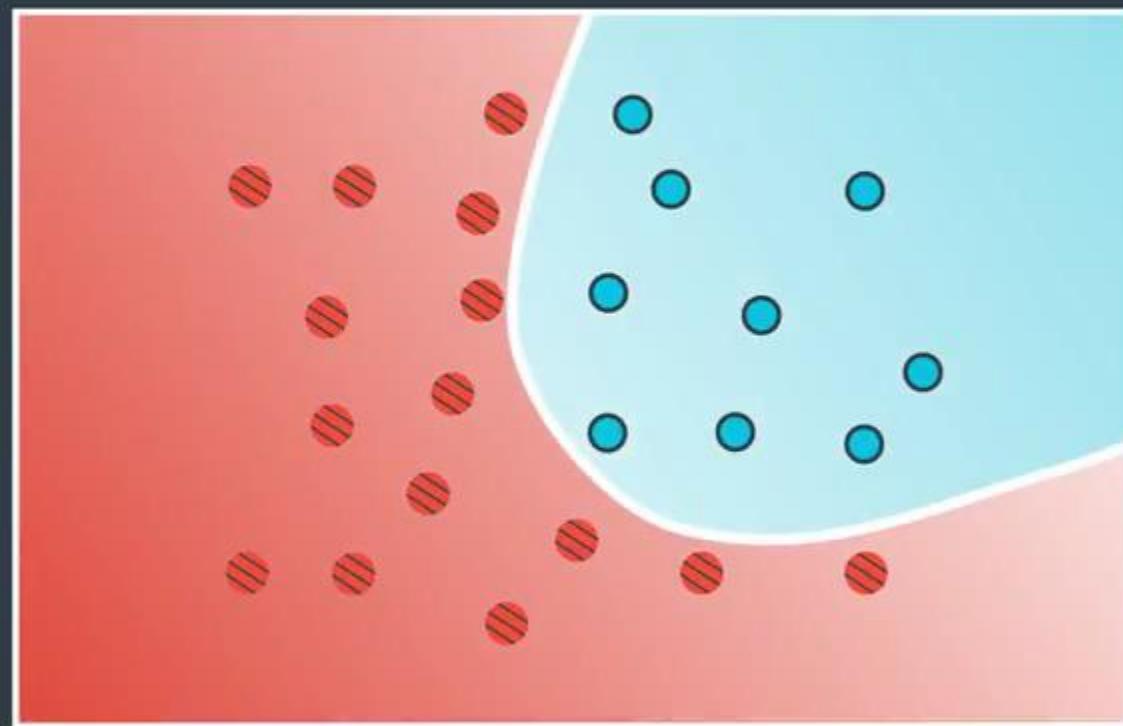
Acceptance at a University



Acceptance at a University

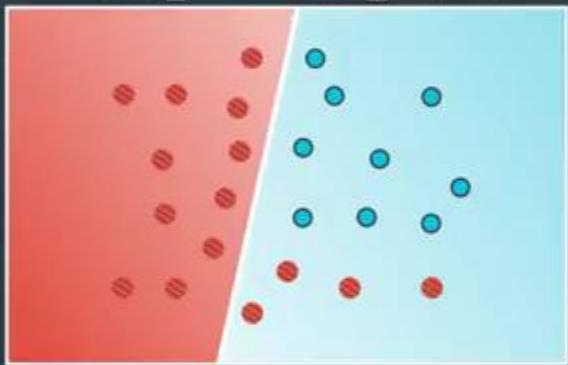


Non-Linear Regions

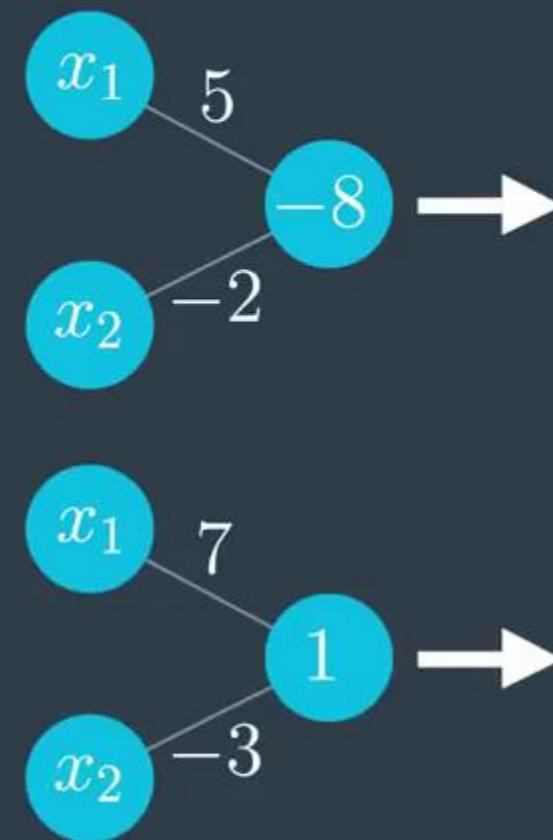
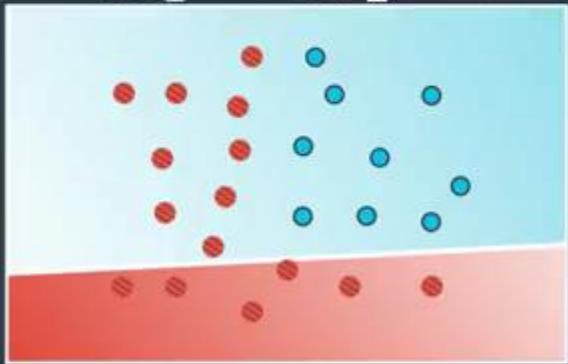


Neural Network

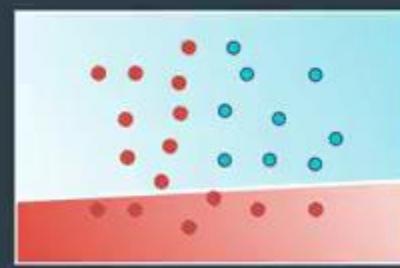
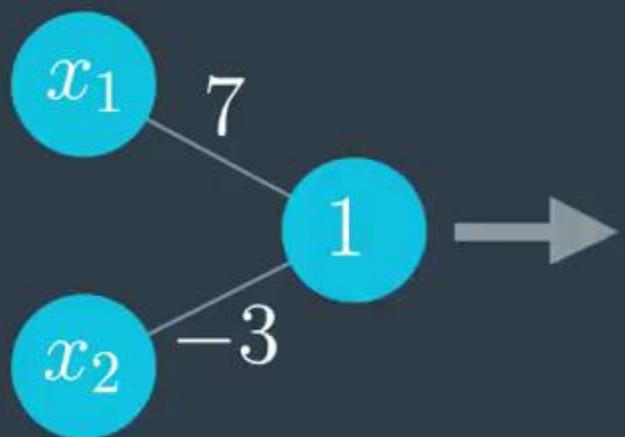
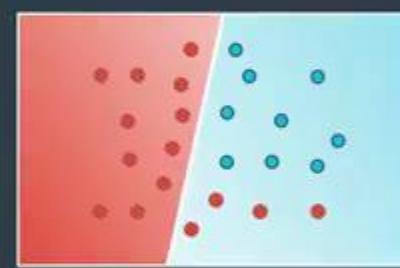
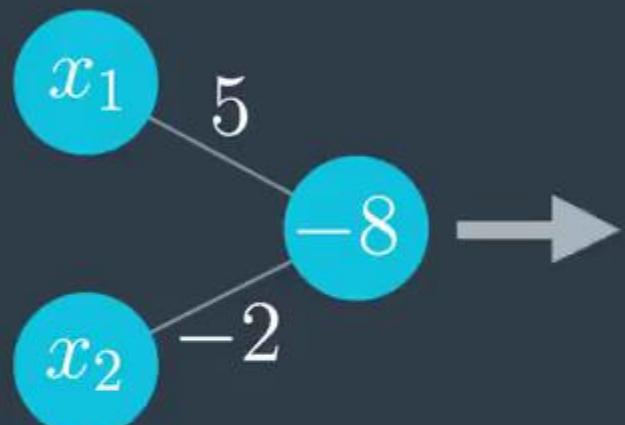
$$5x_1 - 2x_2 + 8$$



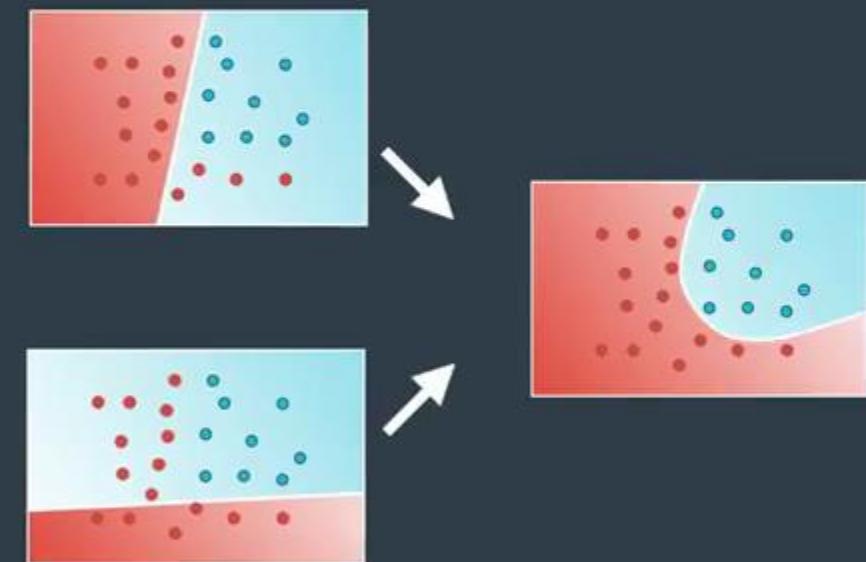
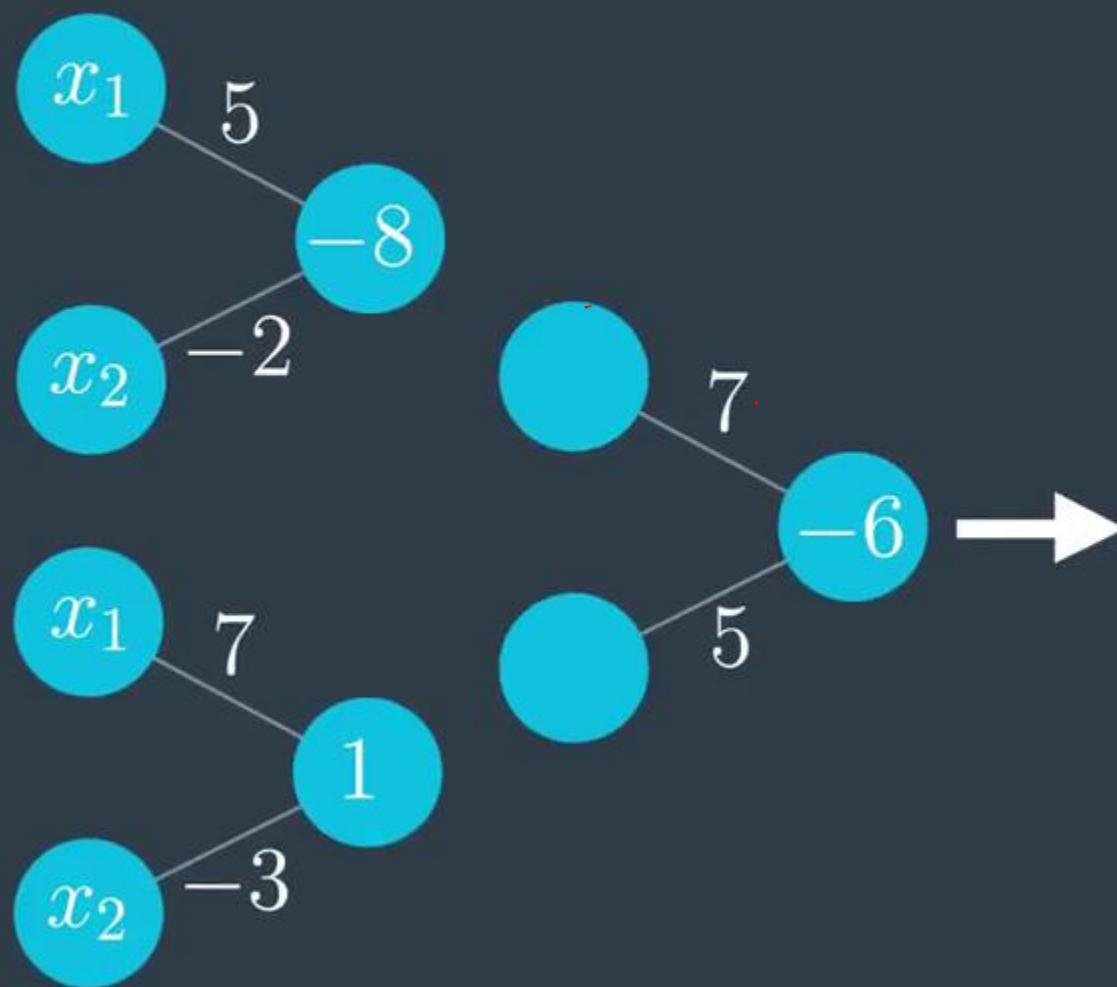
$$7x_1 - 3x_2 - 1$$



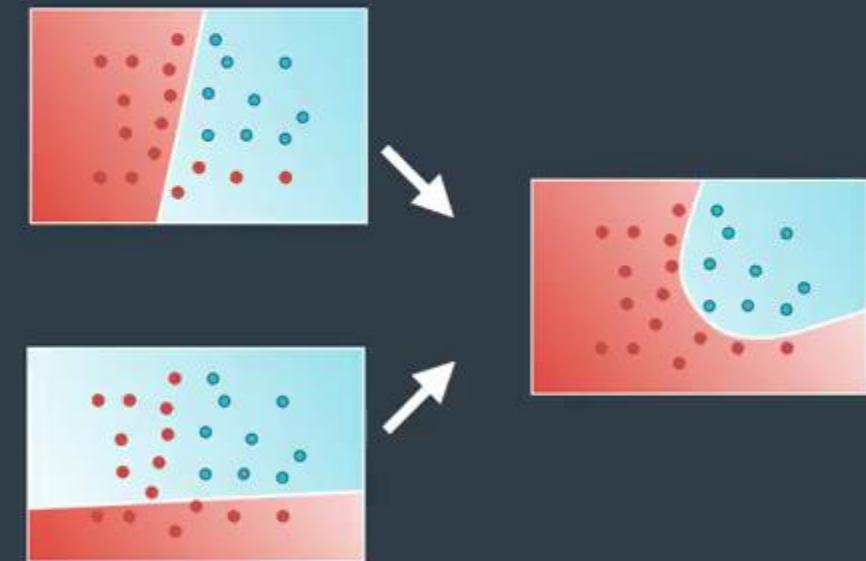
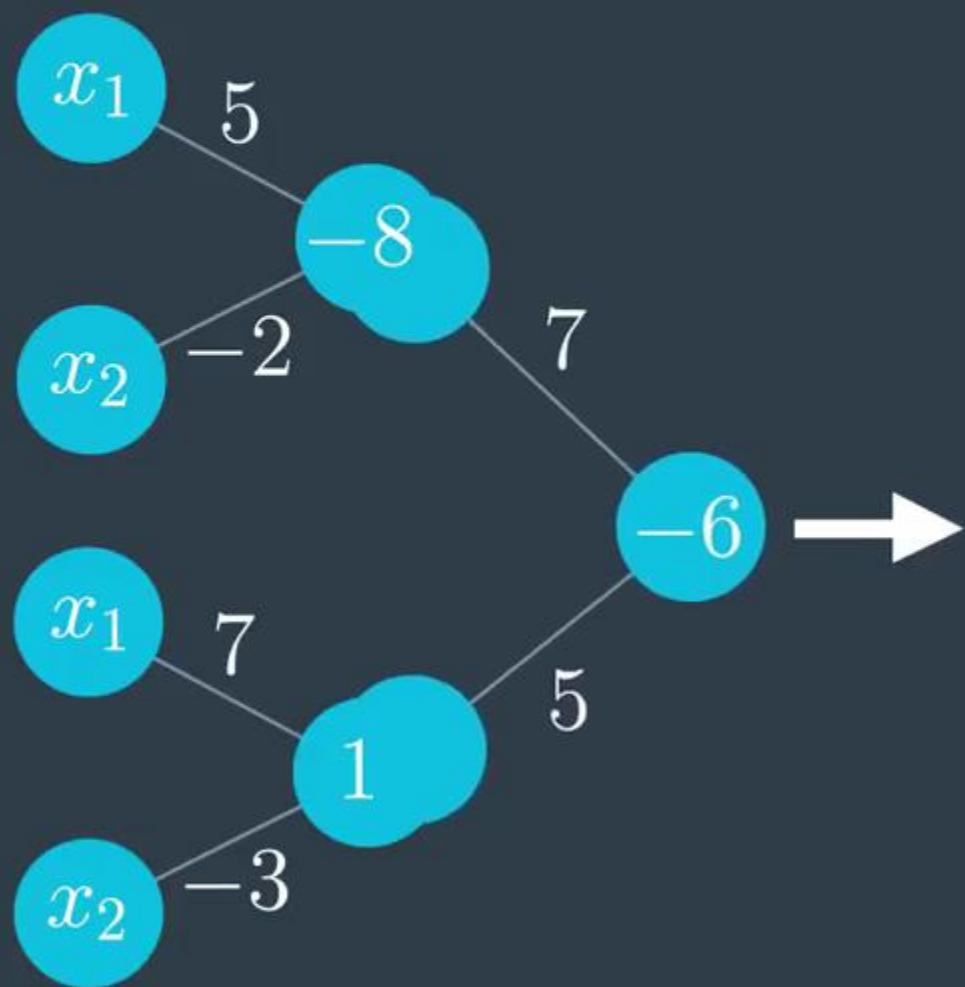
Neural Network



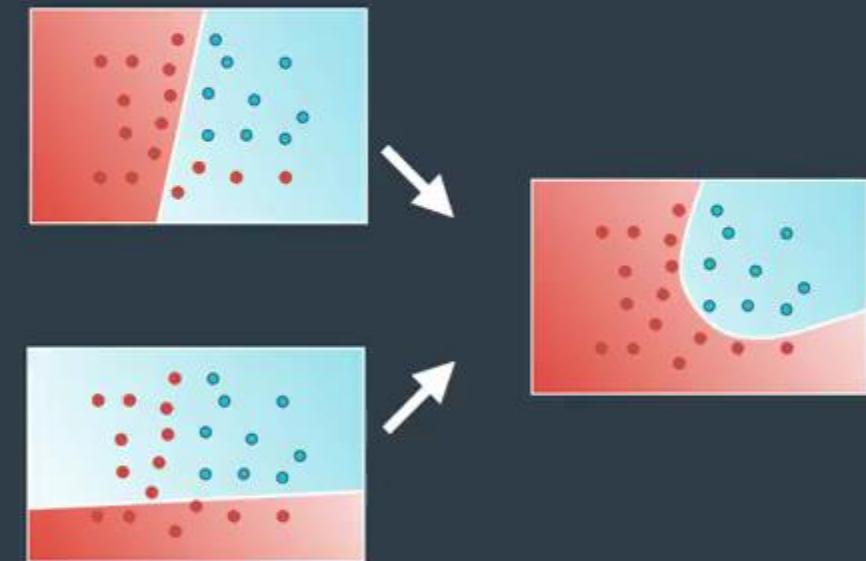
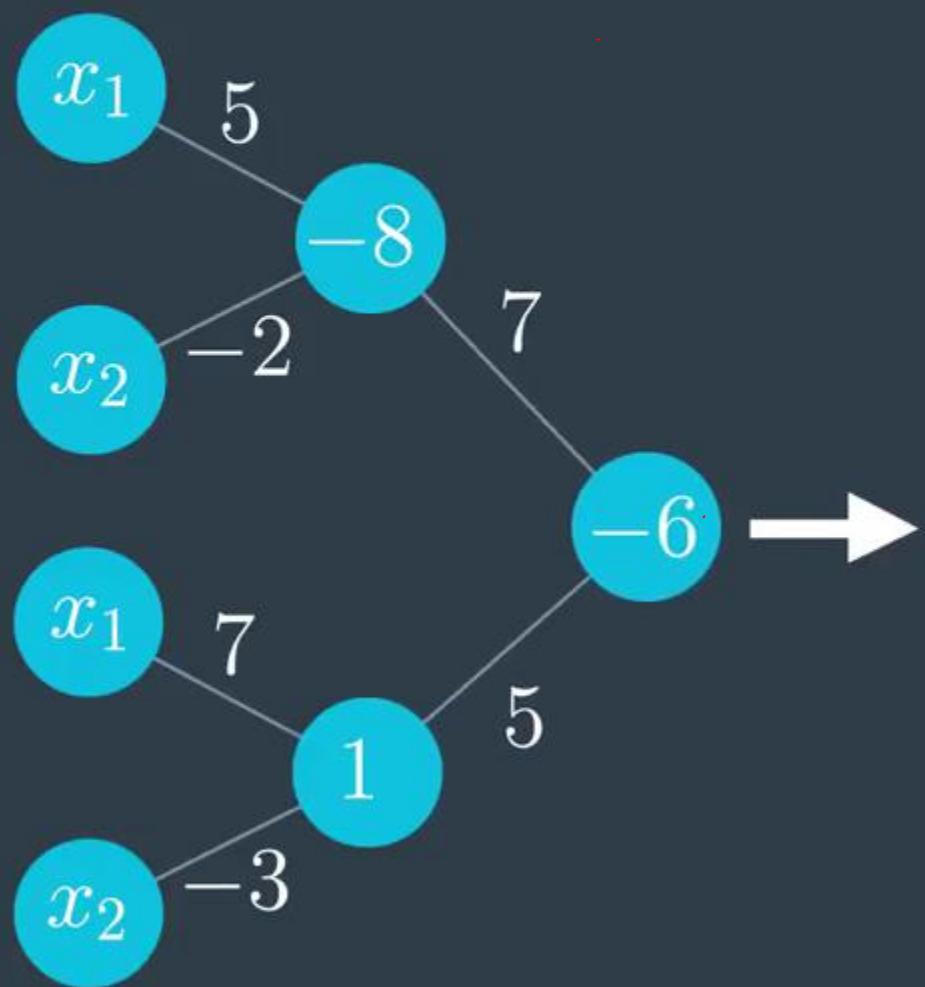
Neural Network



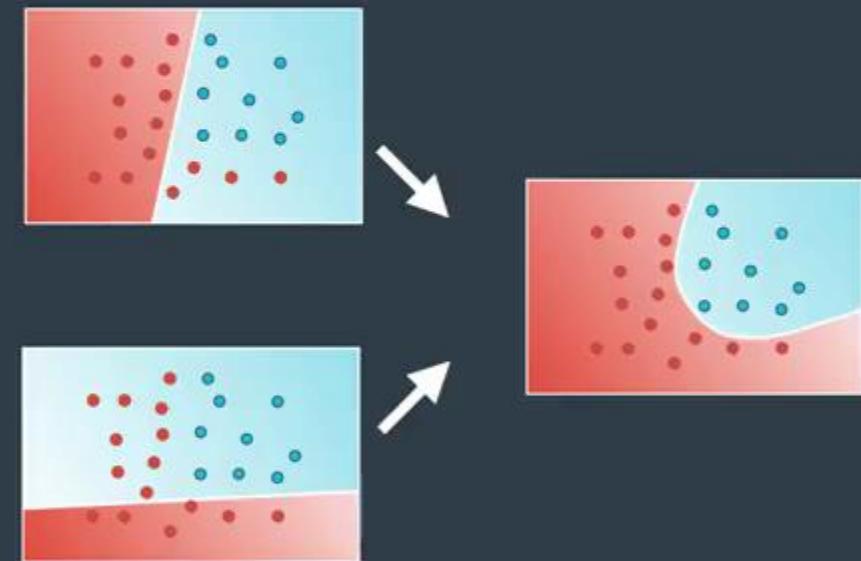
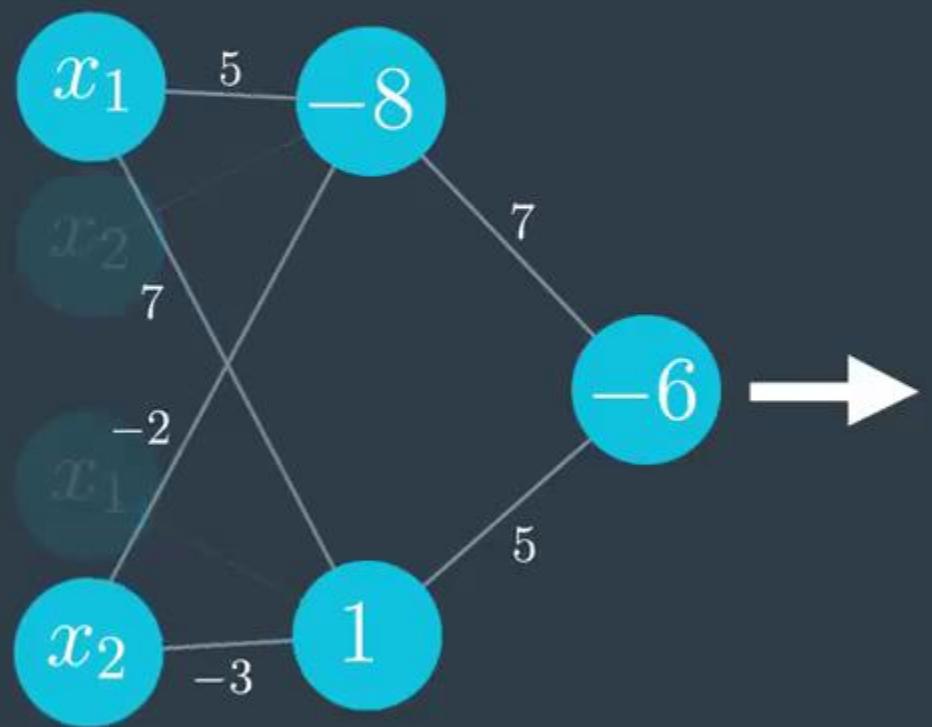
Neural Network



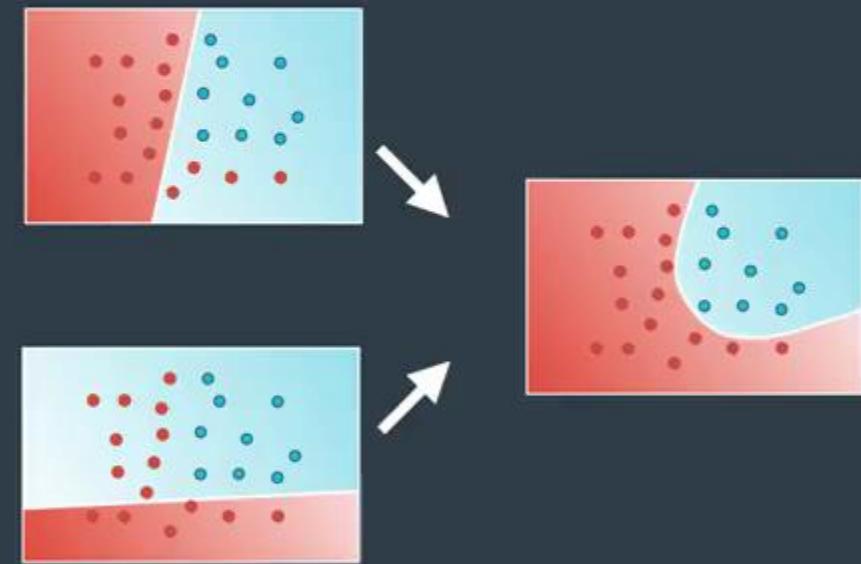
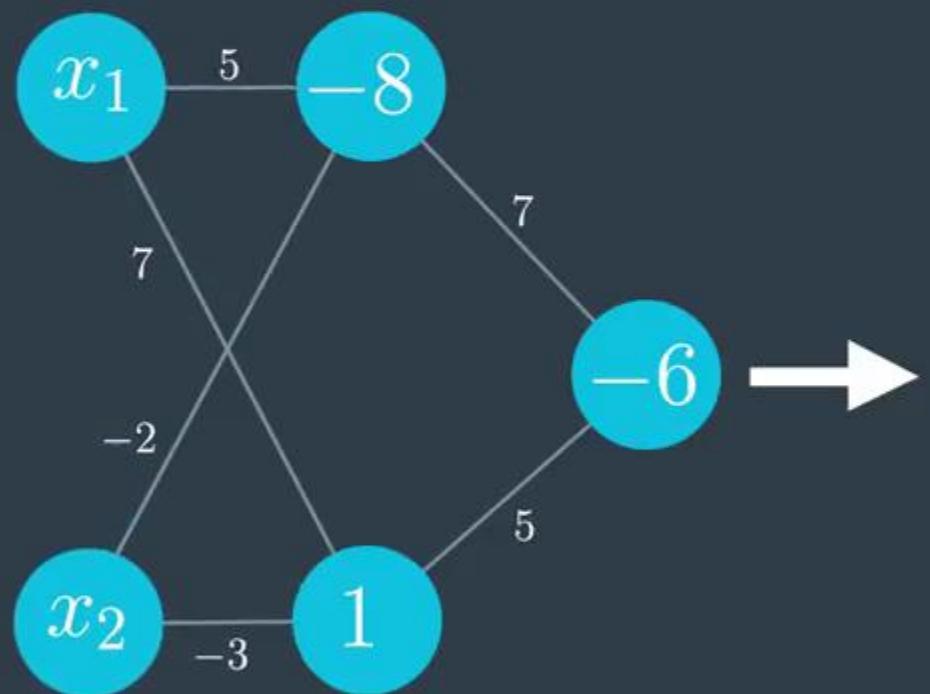
Neural Network



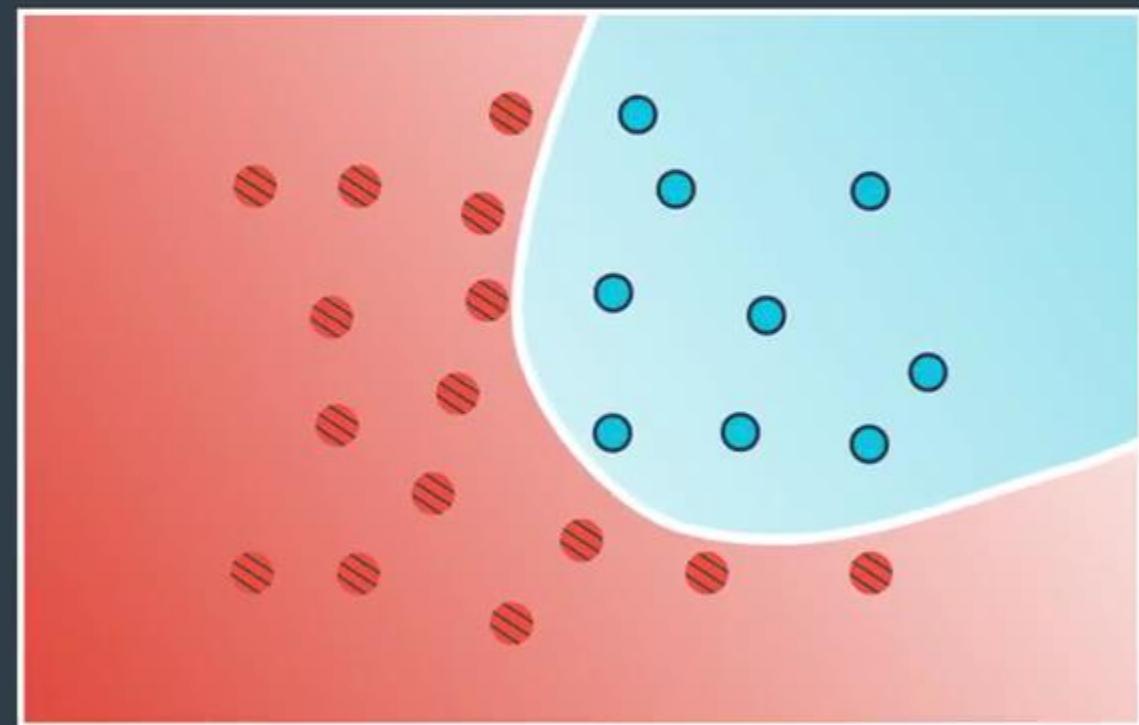
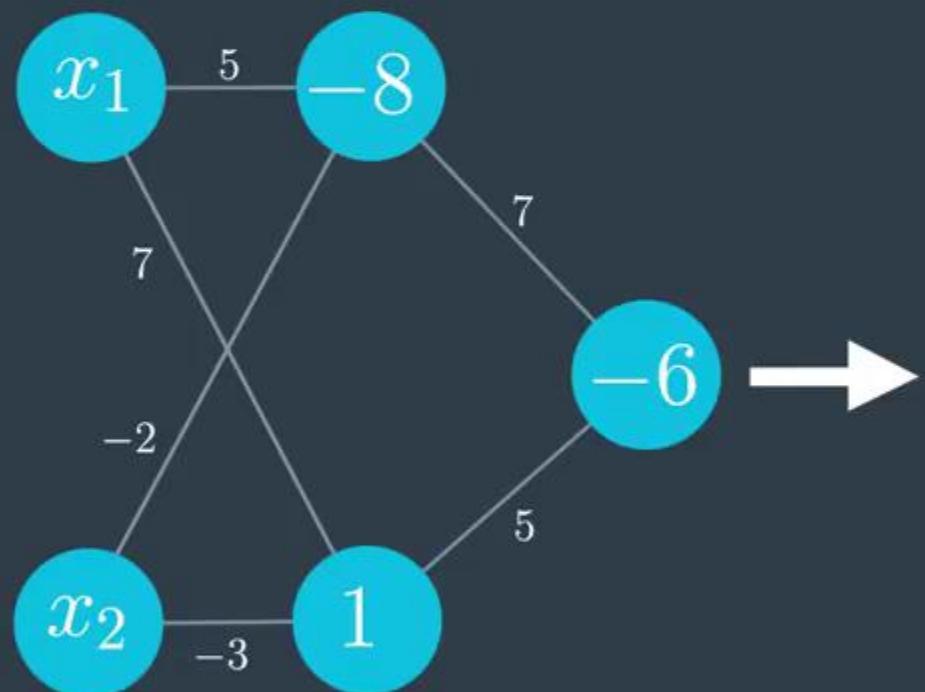
Neural Network



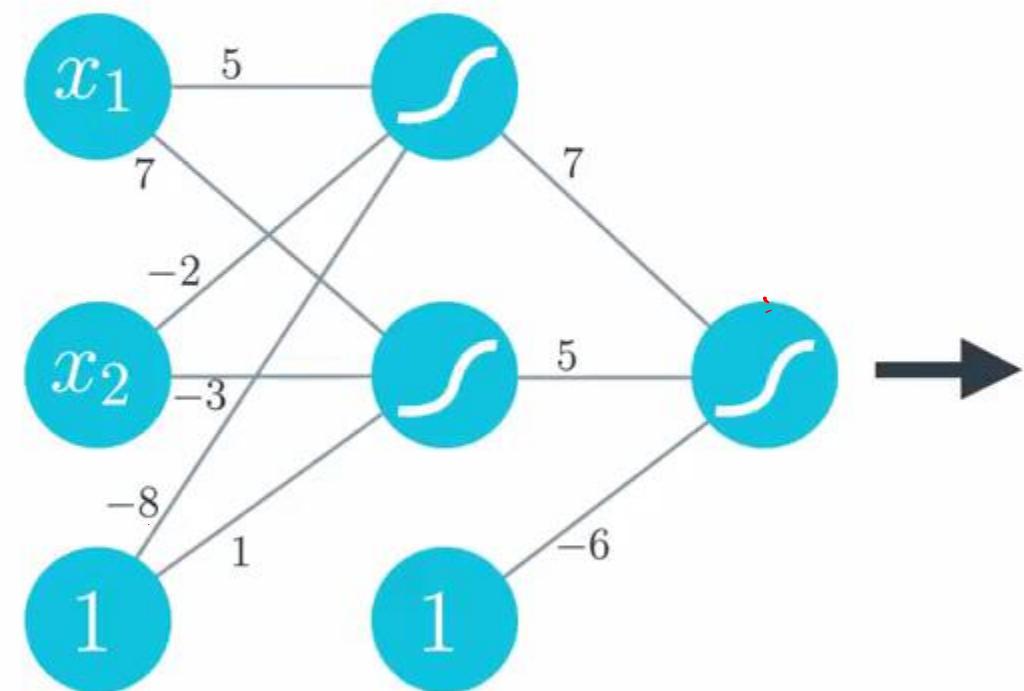
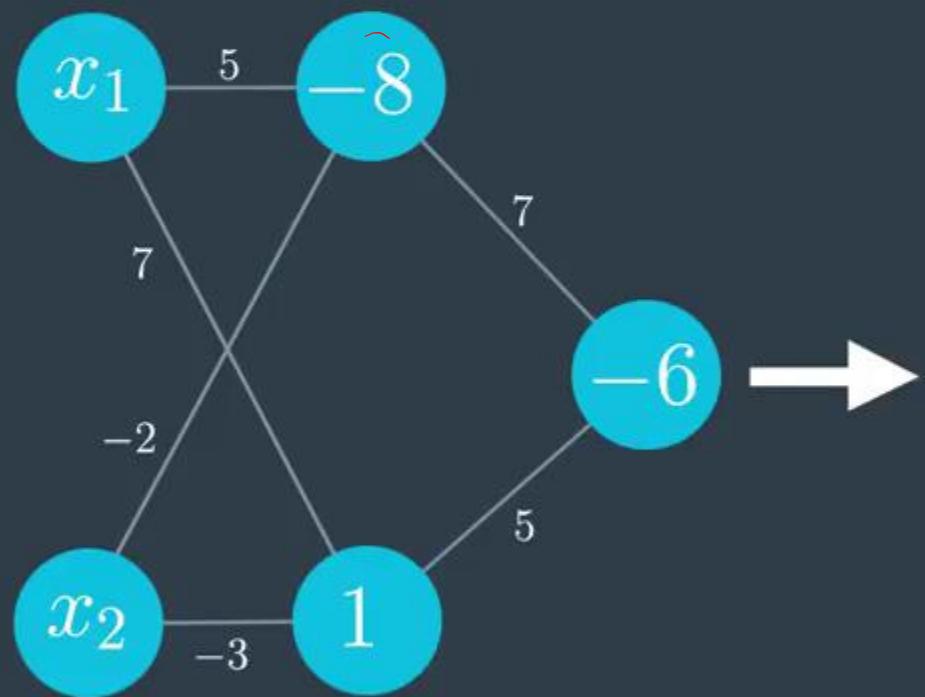
Neural Network



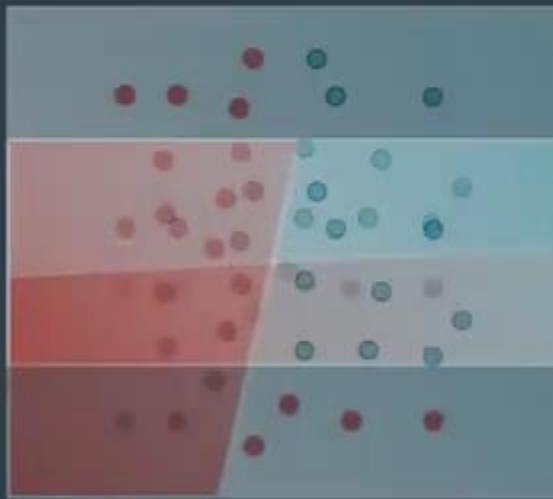
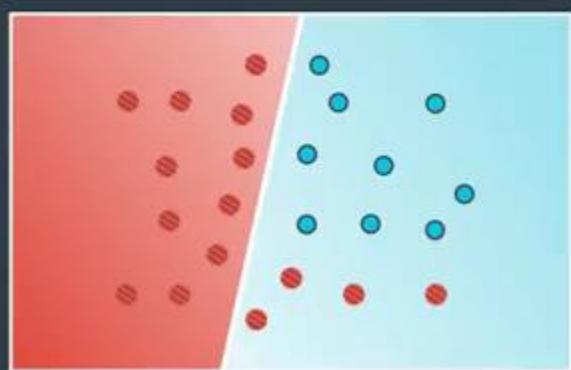
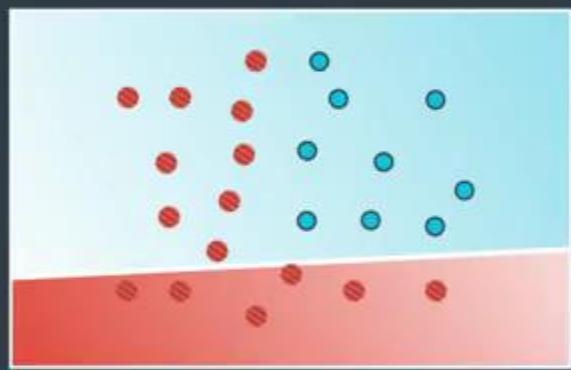
Neural Network



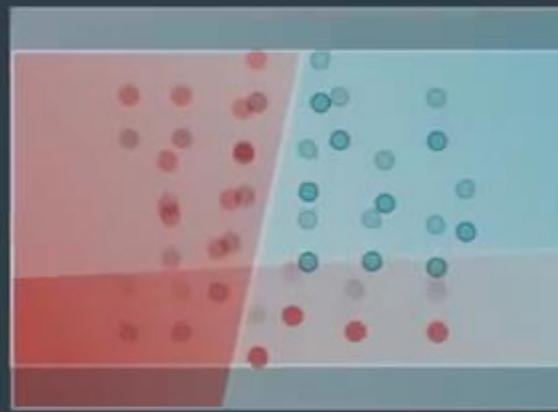
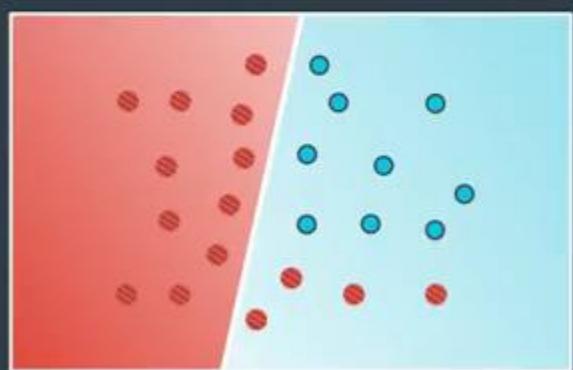
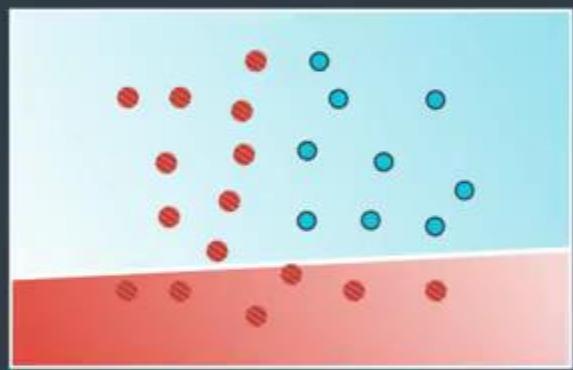
Neural Network



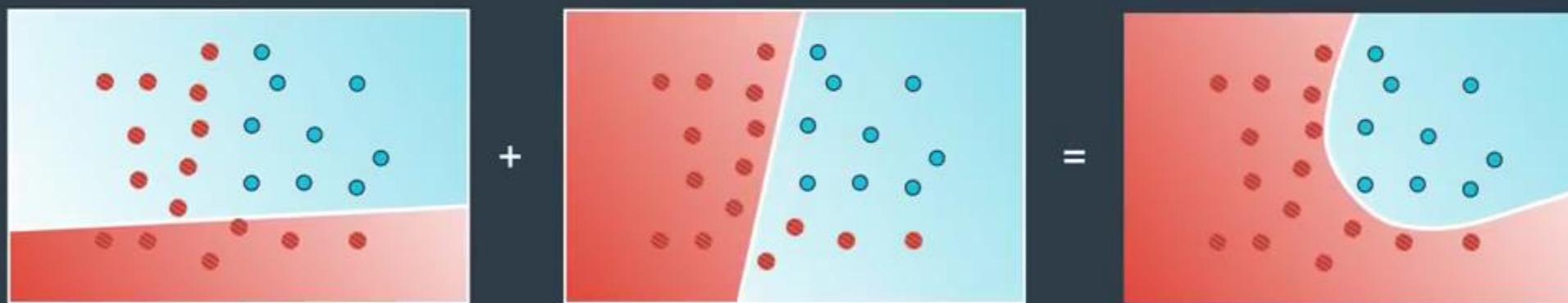
Combining Regions



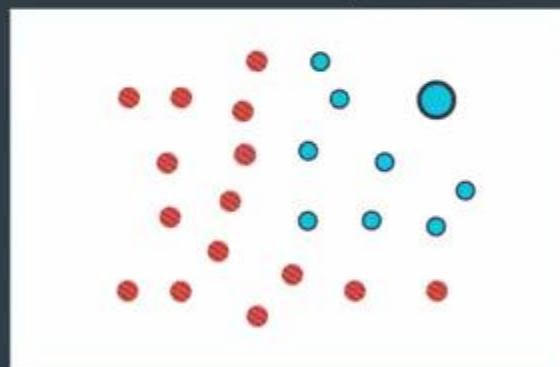
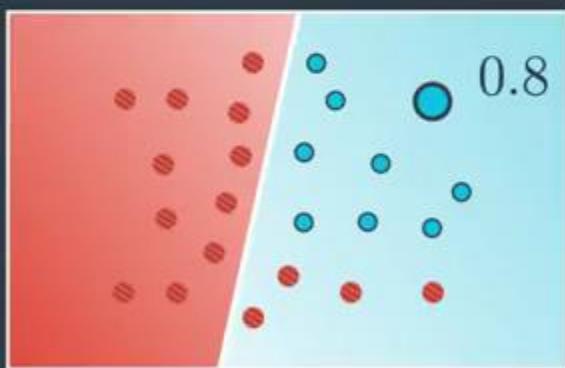
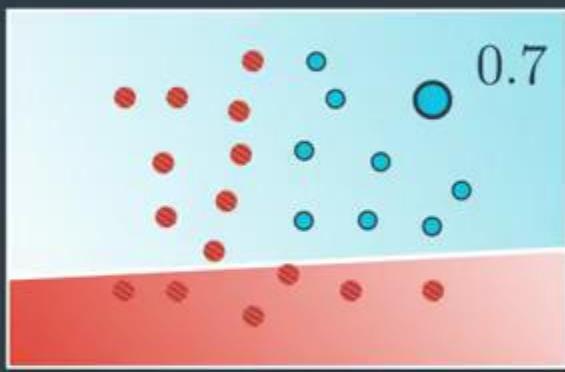
Combining Regions



Combining Regions

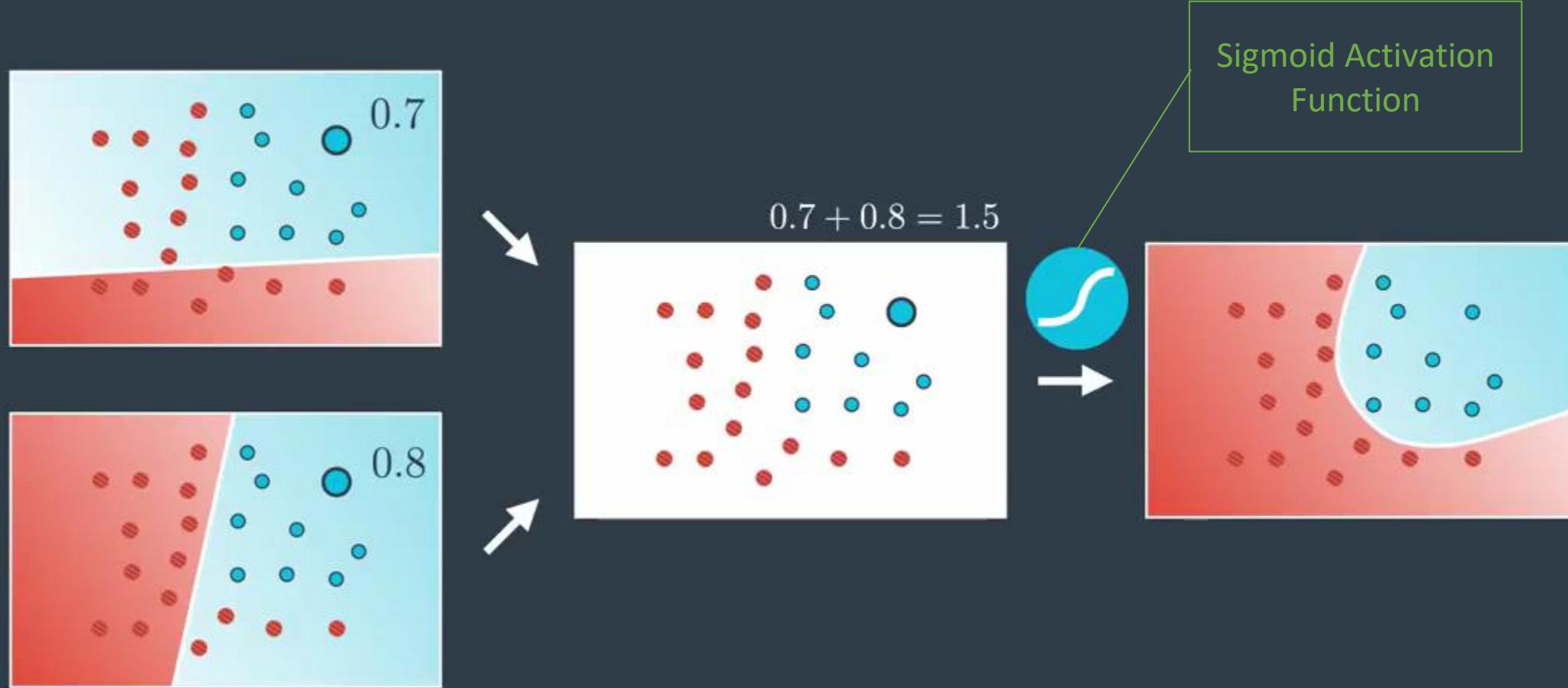


Neural Network

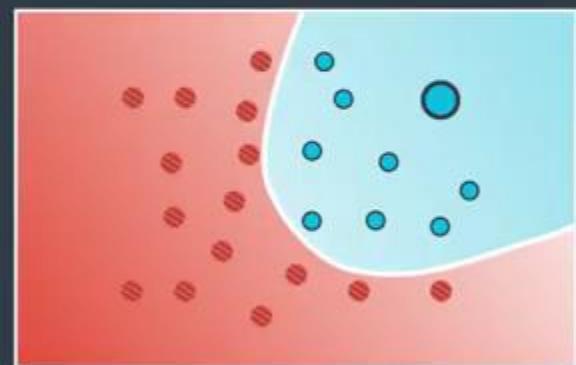
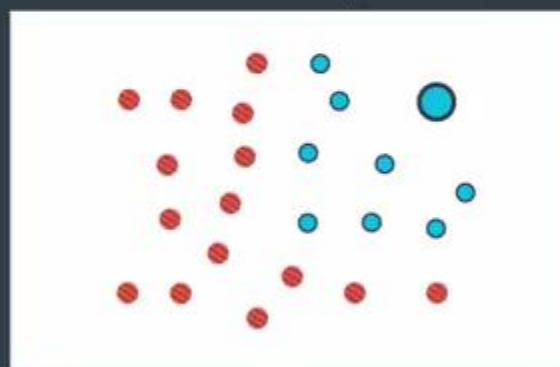
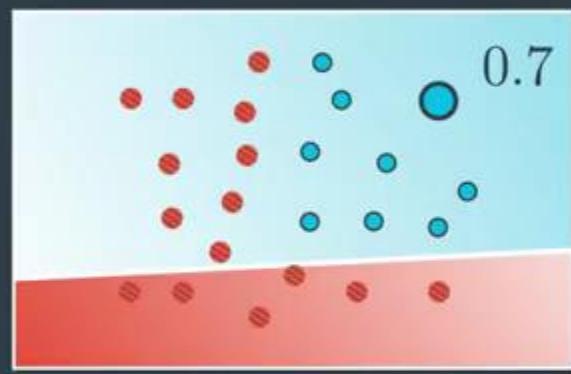


It's bigger than 1
We have to do
something

Neural Network

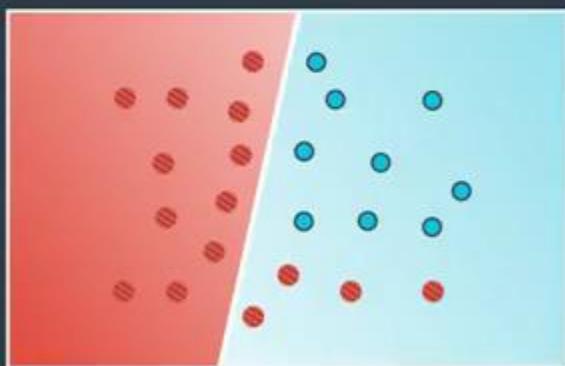
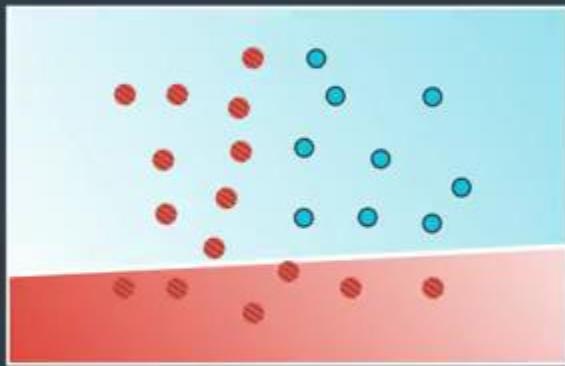


Neural Network



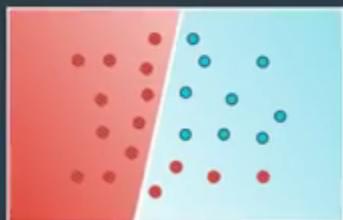
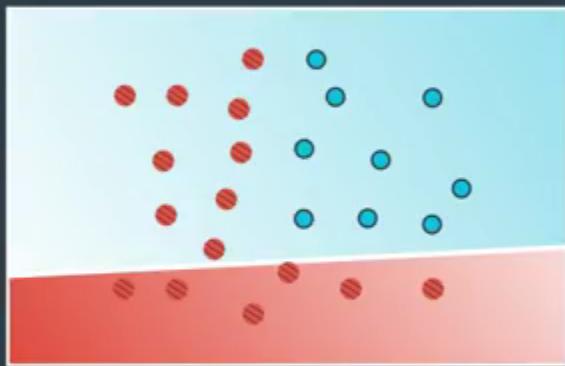
Combining Regions

What Weights can do?



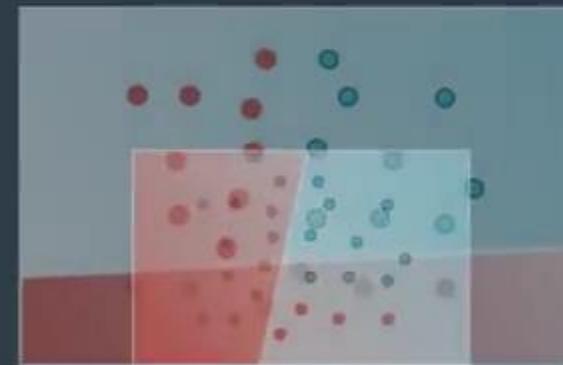
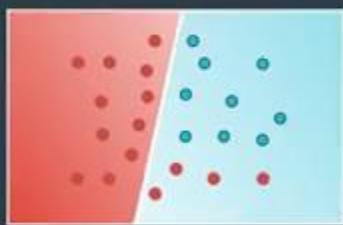
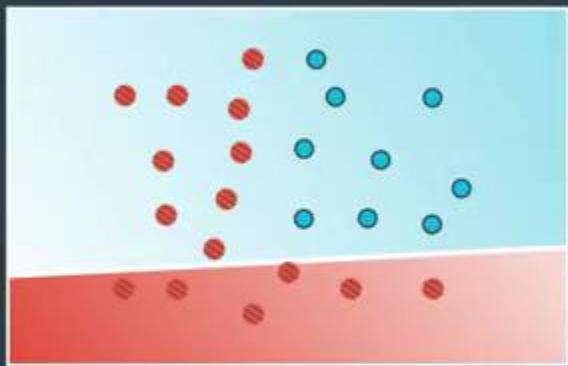
Combining Regions

What Weights can do?



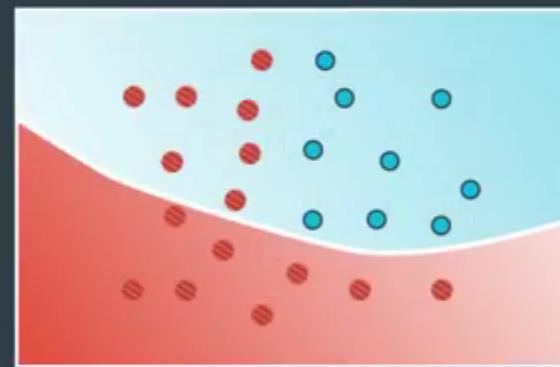
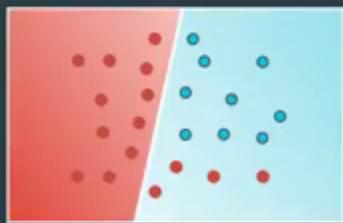
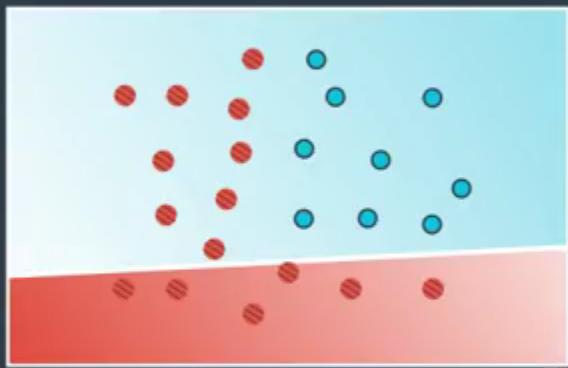
Combining Regions

What Weights can do?



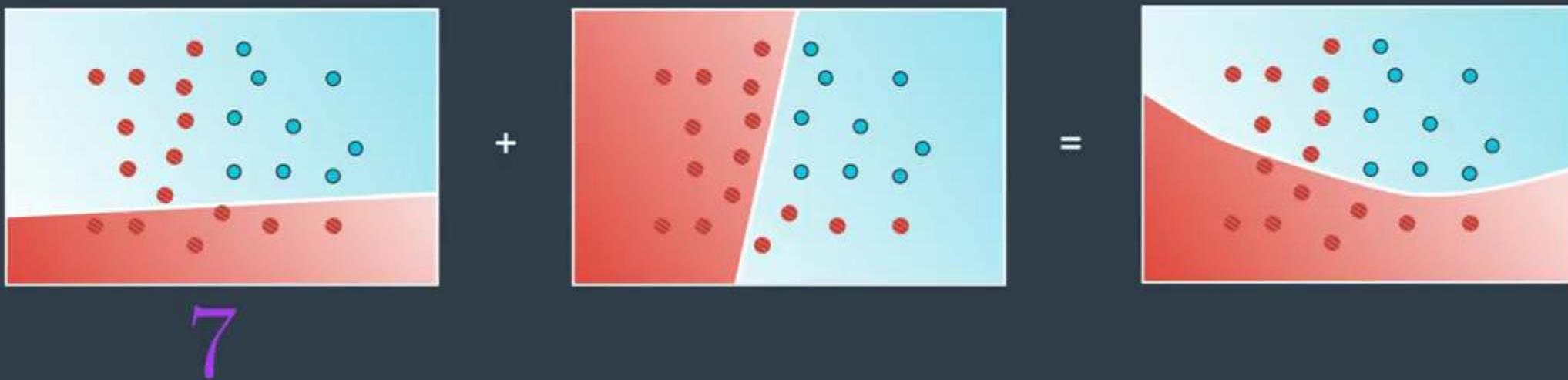
Combining Regions

What Weights can do?



Combining Regions

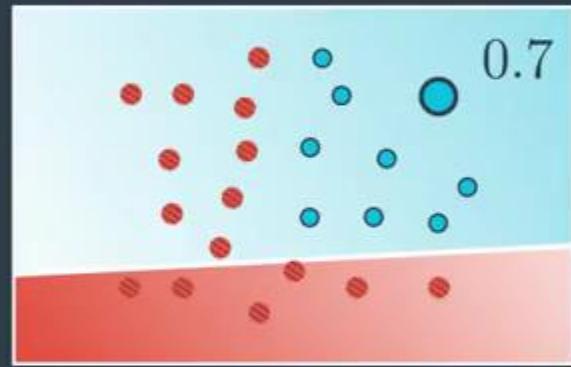
What Weights can do?



Neural Network

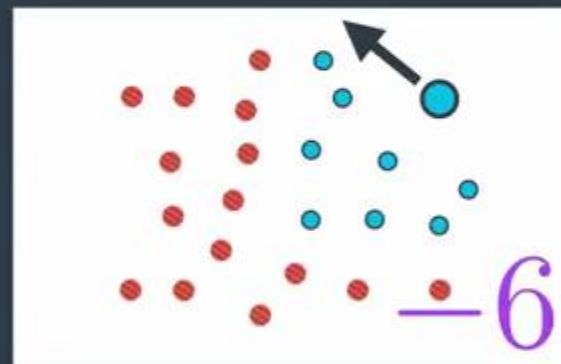
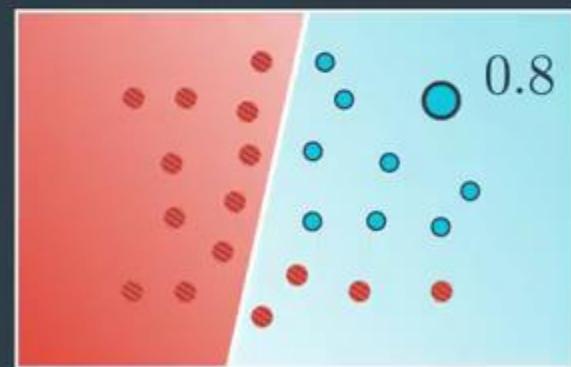
What Weights can do?

7



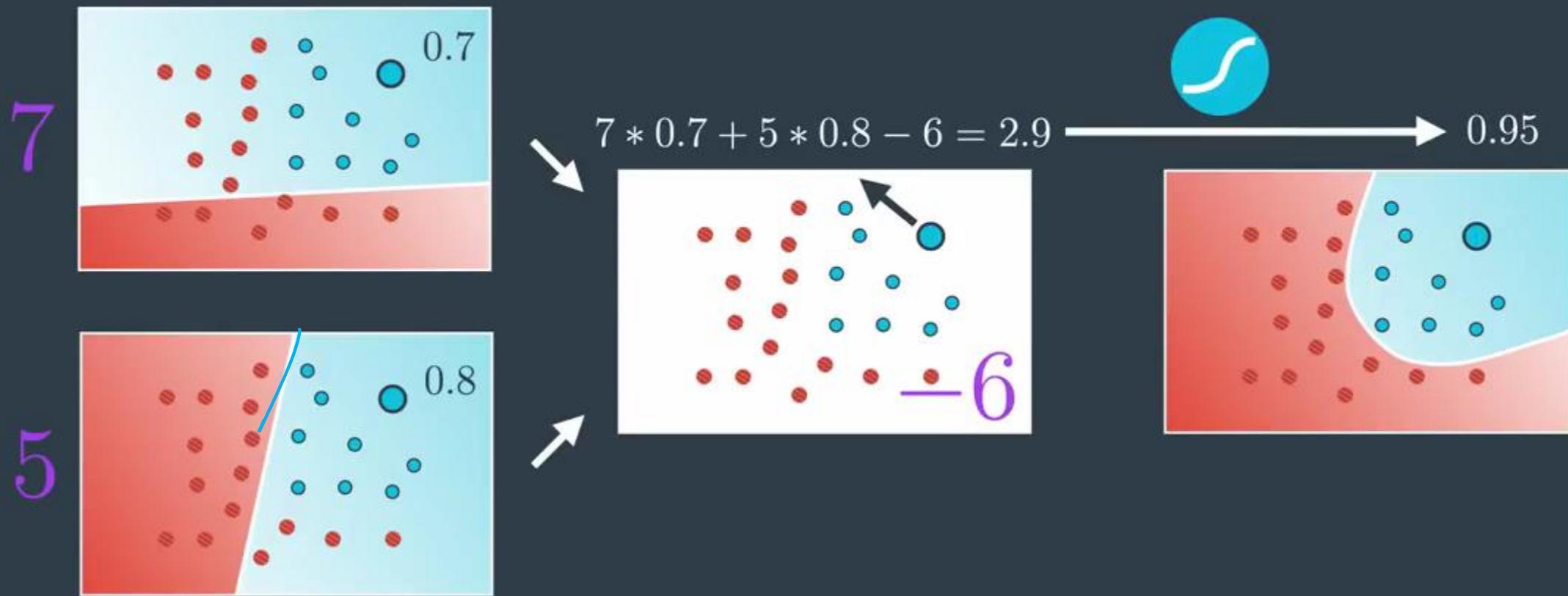
$$7 * 0.7 + 5 * 0.8 - 6 = 2.9$$

5



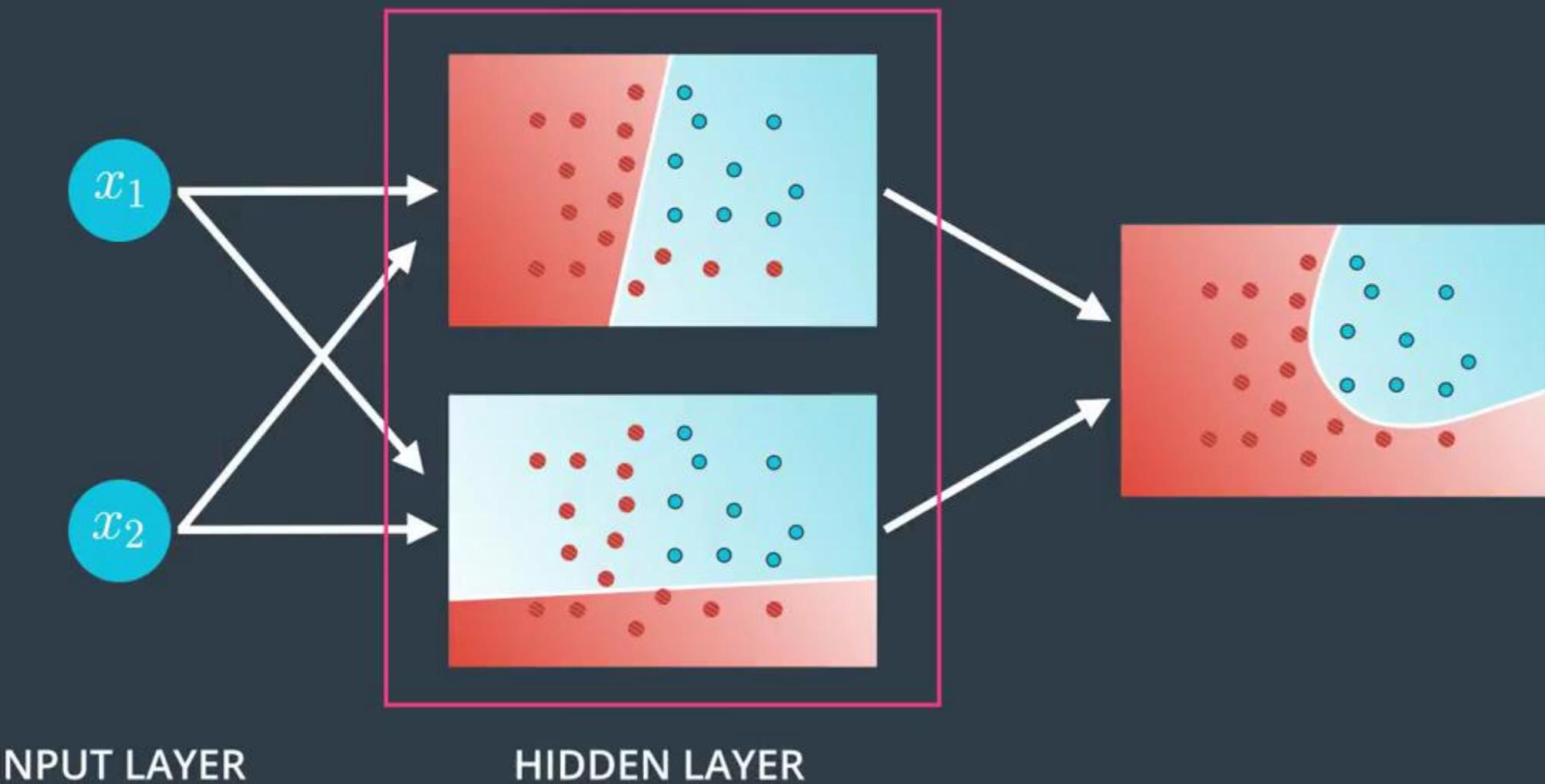
Neural Network

What Weights can do?



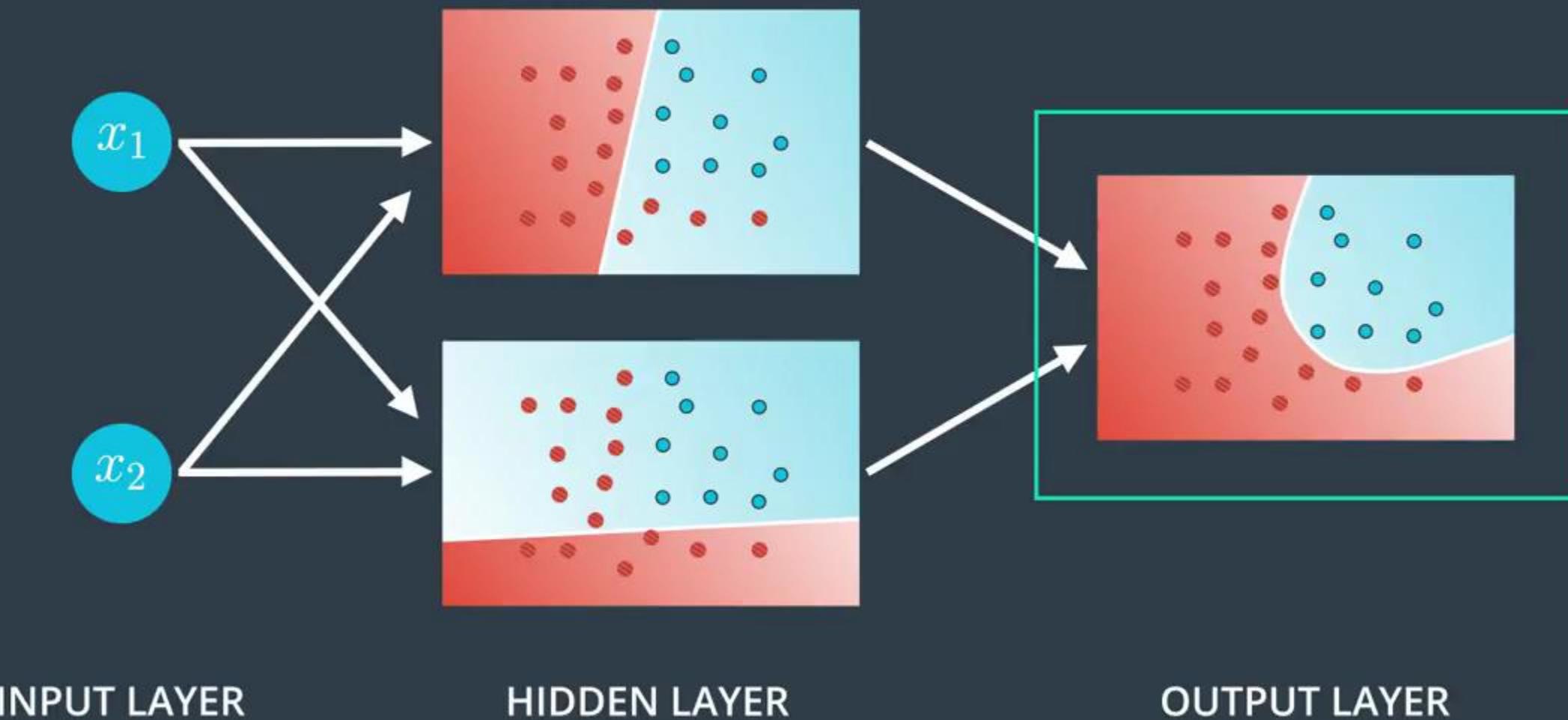
What Layers can do?

Neural Network

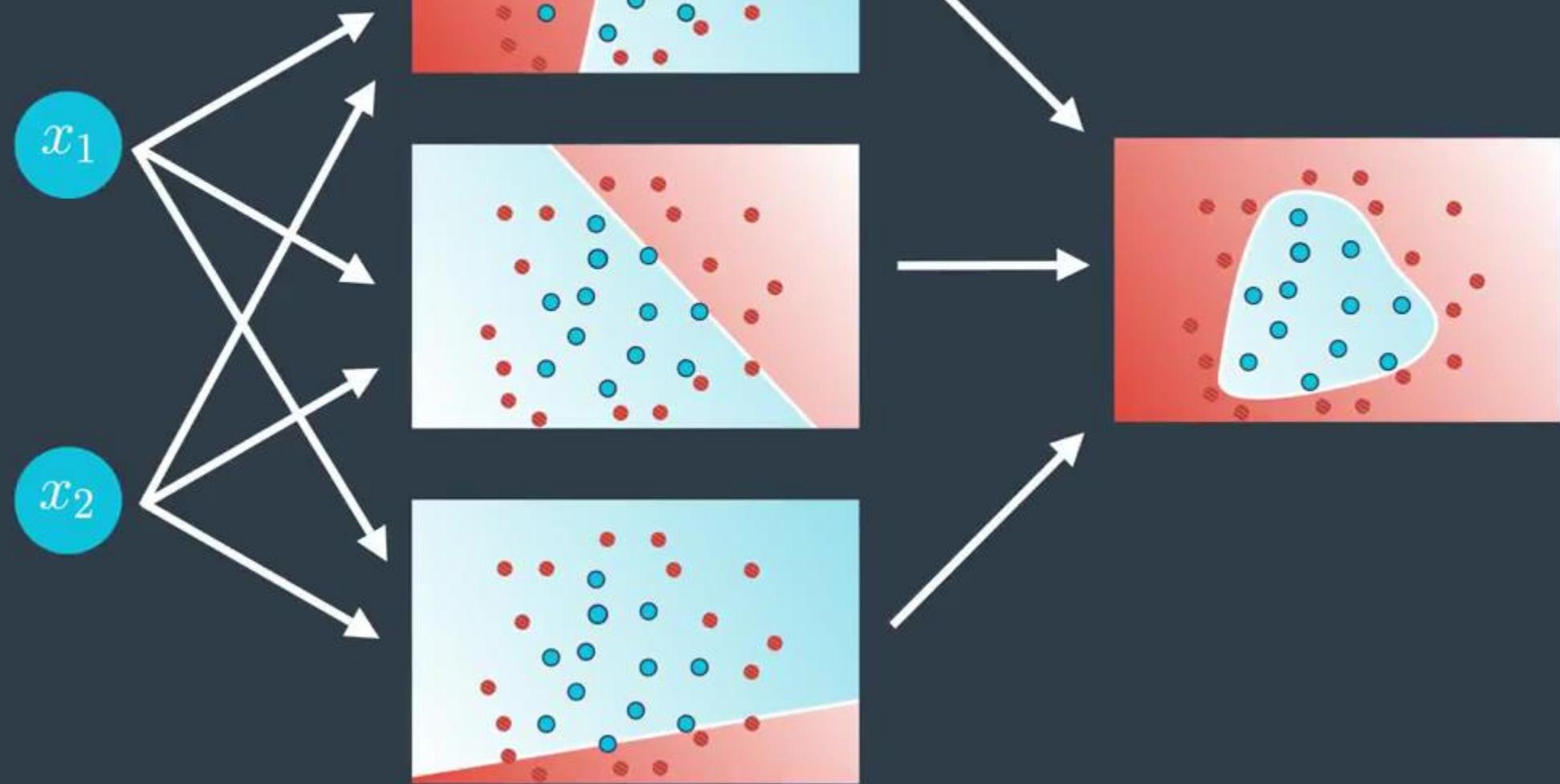


What Layers can do?

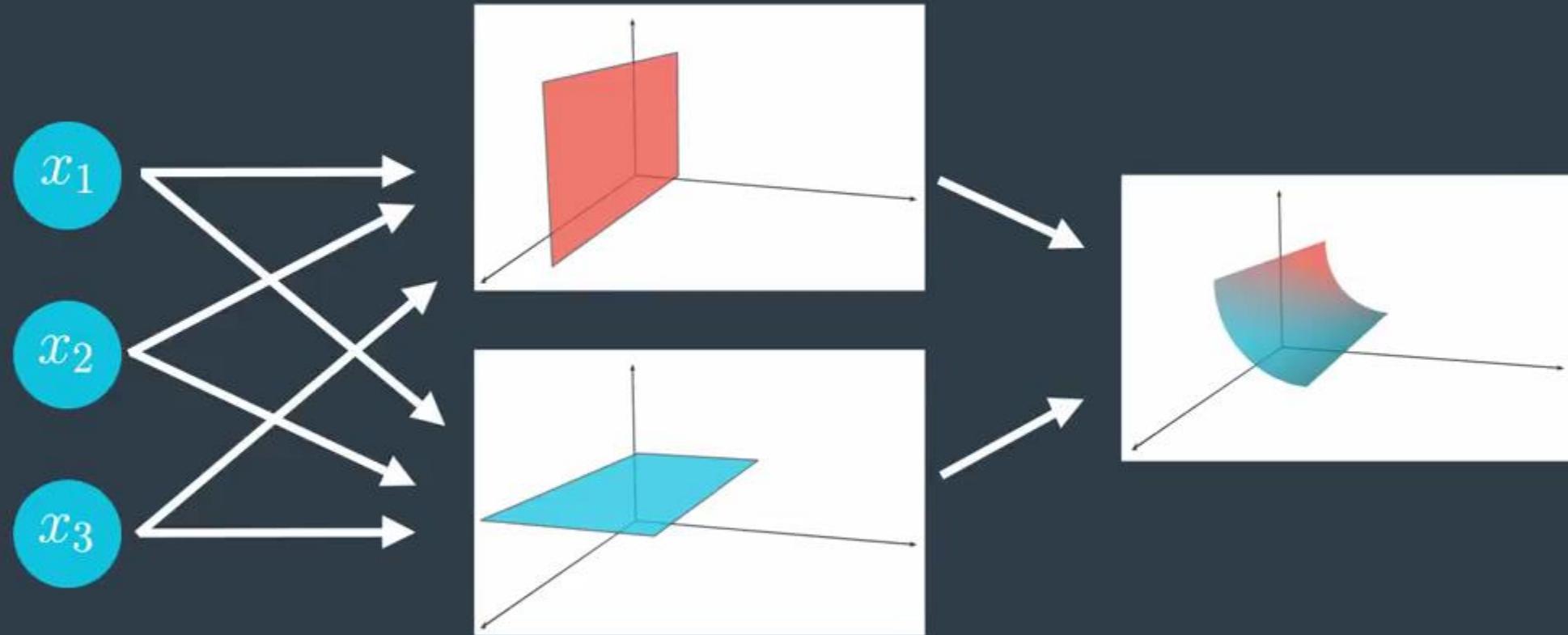
Neural Network



What Layers can do?

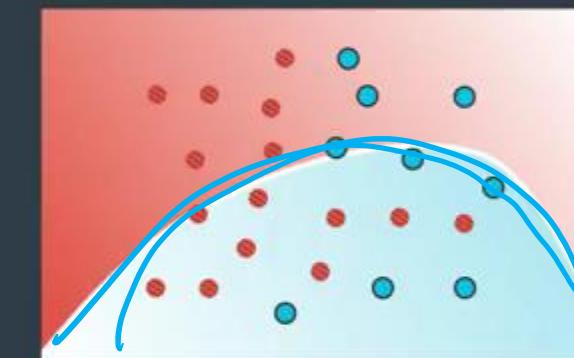
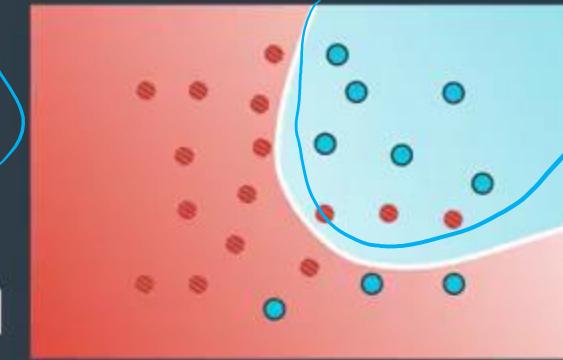
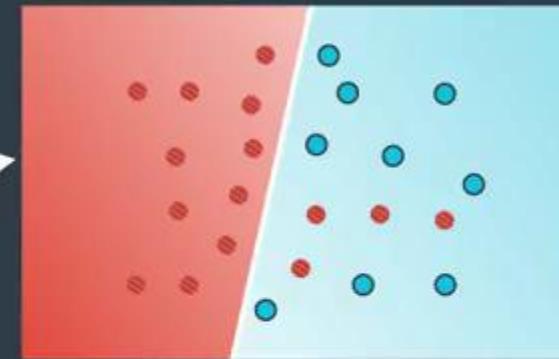
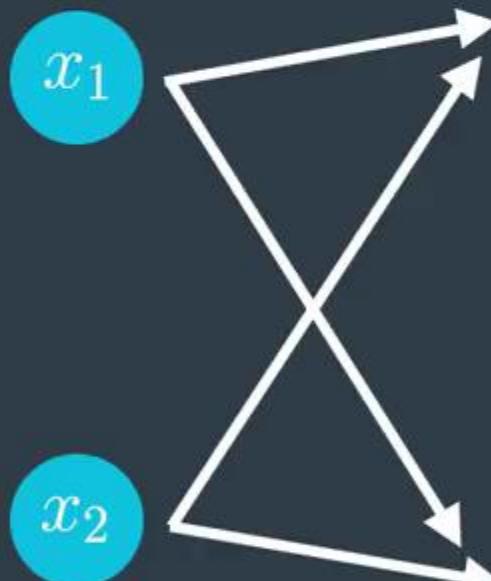


What Layers can do?



Deep Neural Network

Multi-Class Classification



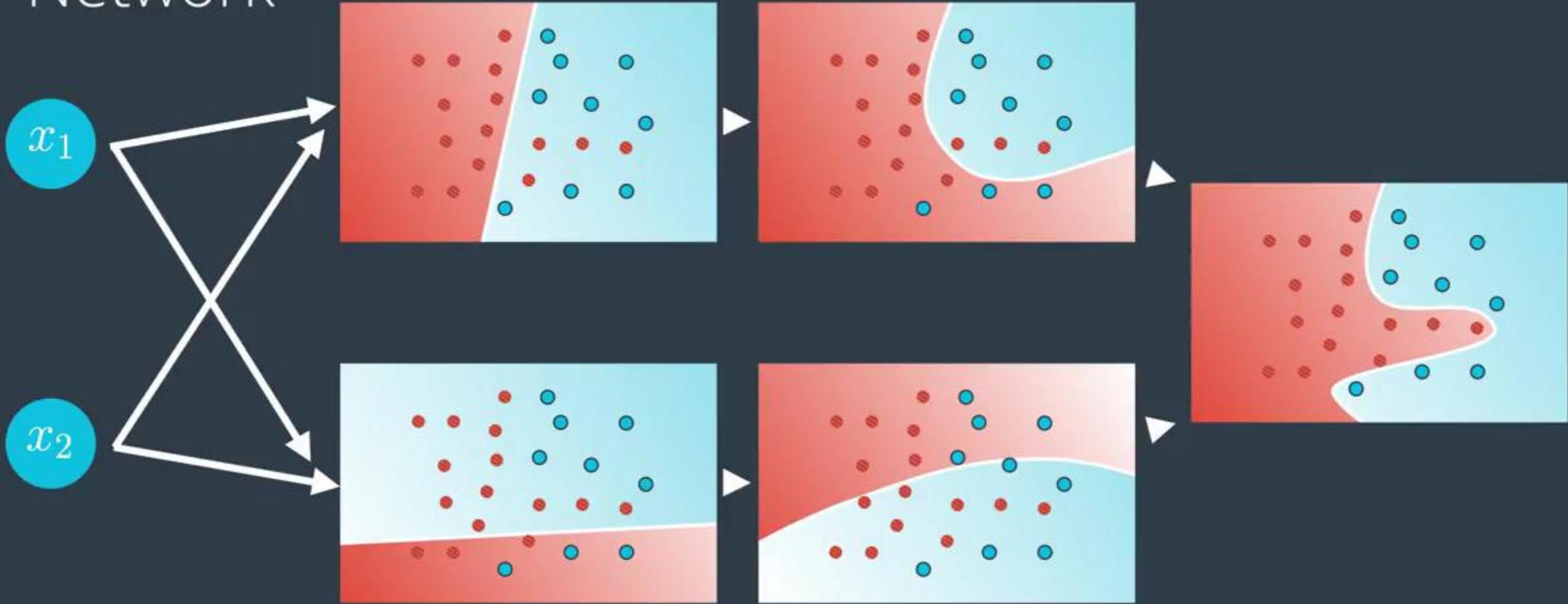
CAT

DOG

BIRD

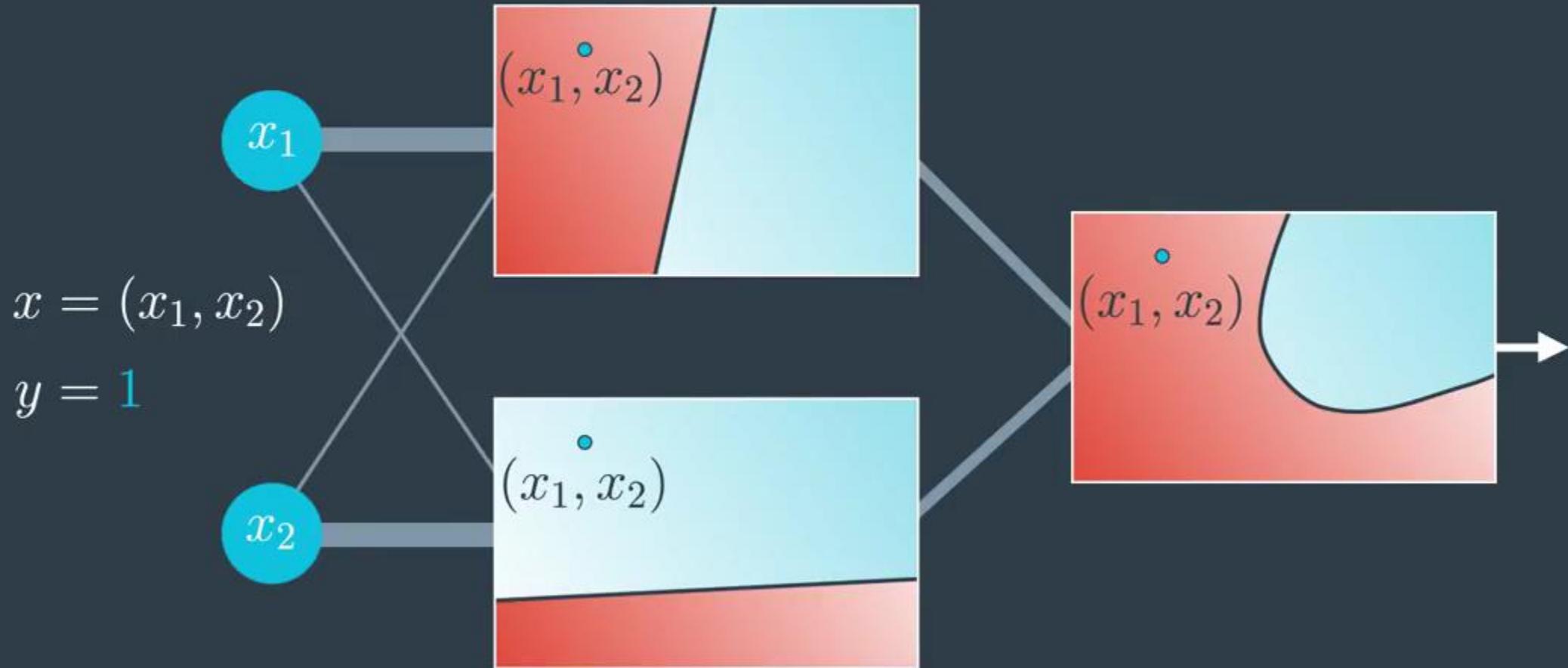
Deep Neural Network

Multi-Layer non-linear
Classification

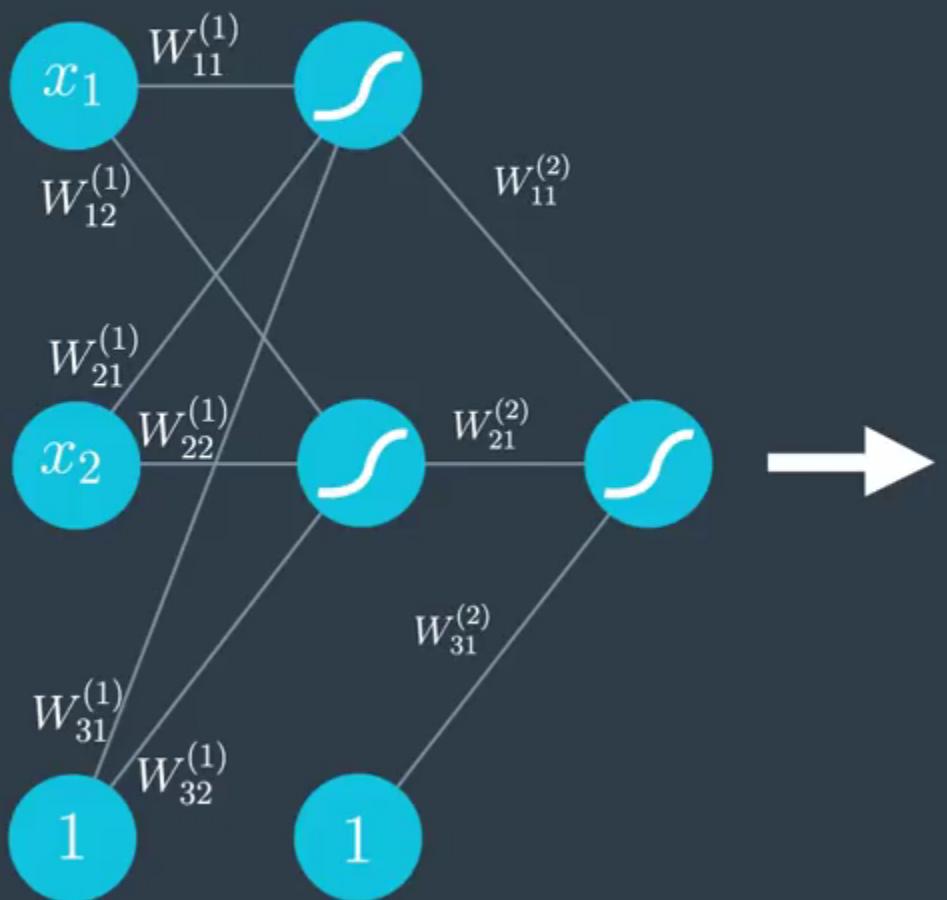


Neural Network

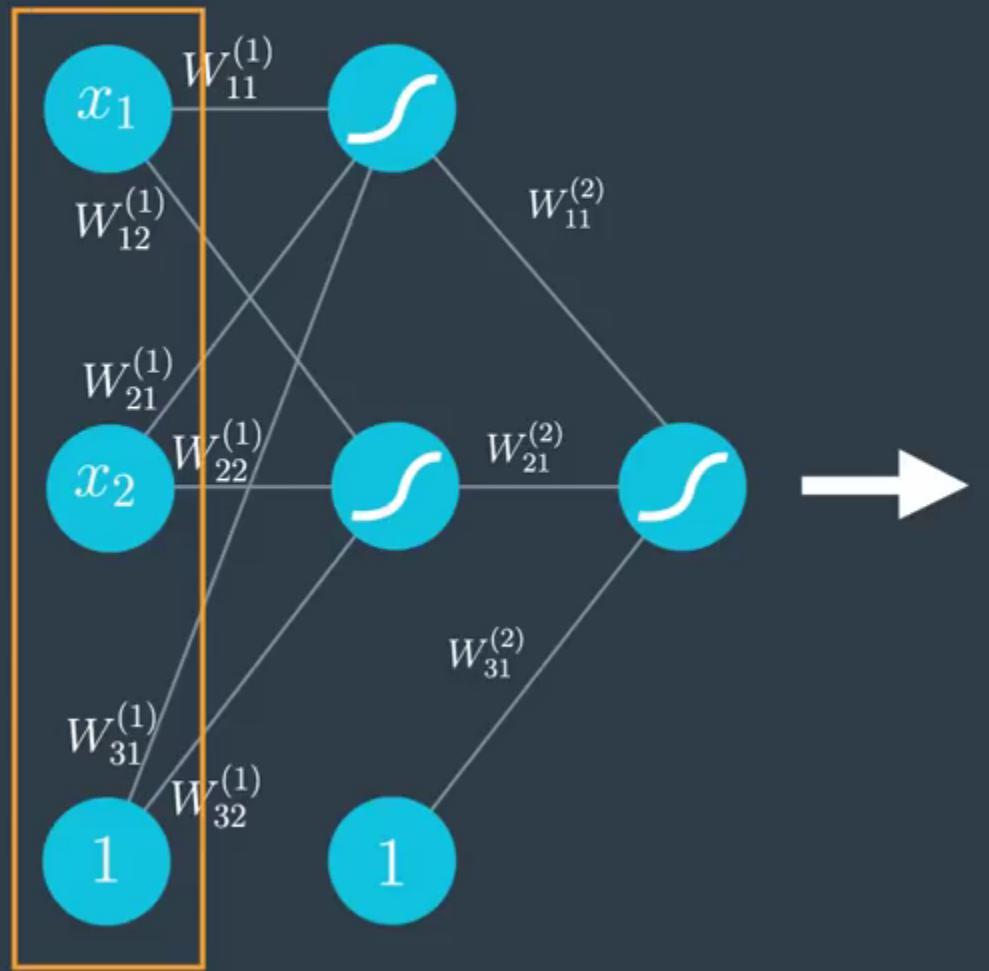
Feed-Forward NN



Feedforward

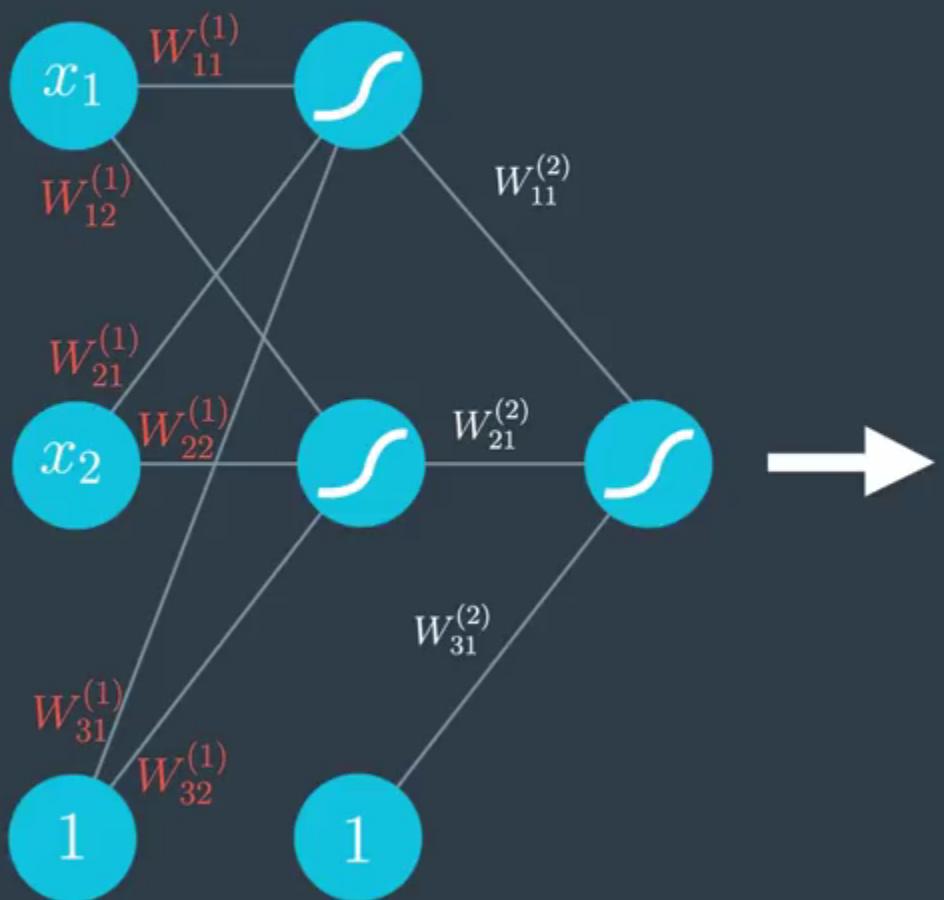


Feedforward



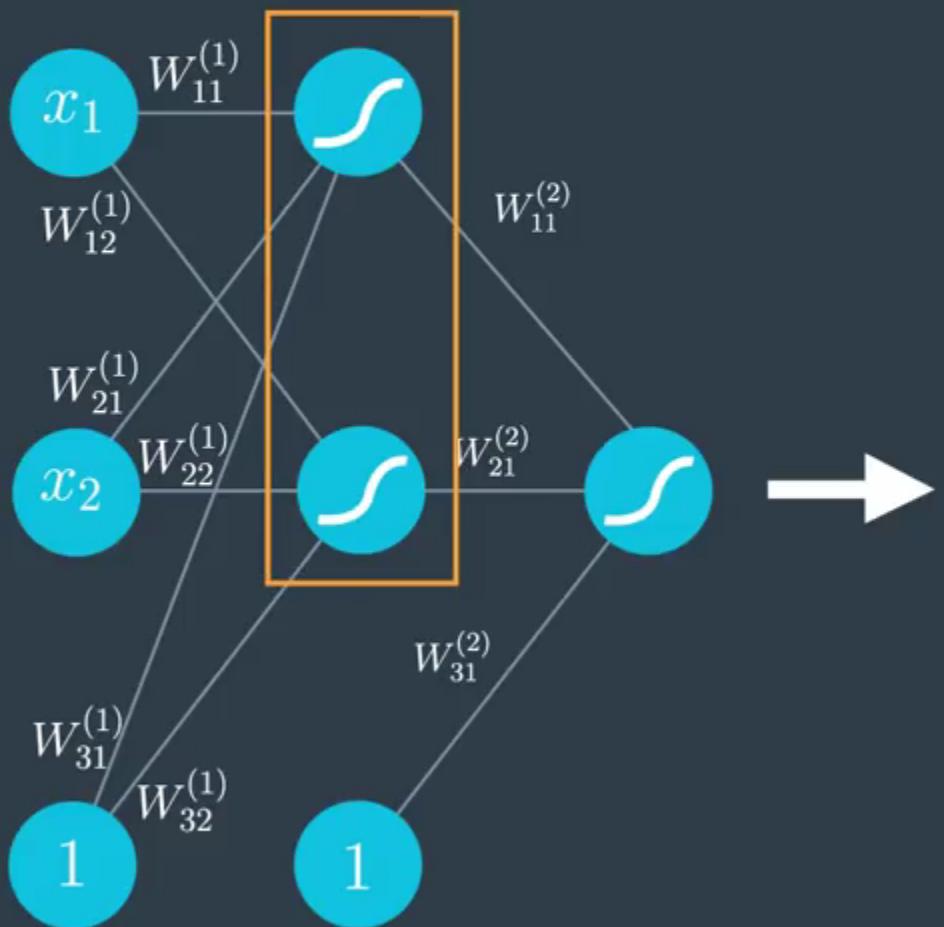
$$\begin{pmatrix} x_1 \\ x_2 \\ 1 \end{pmatrix}$$

Feedforward



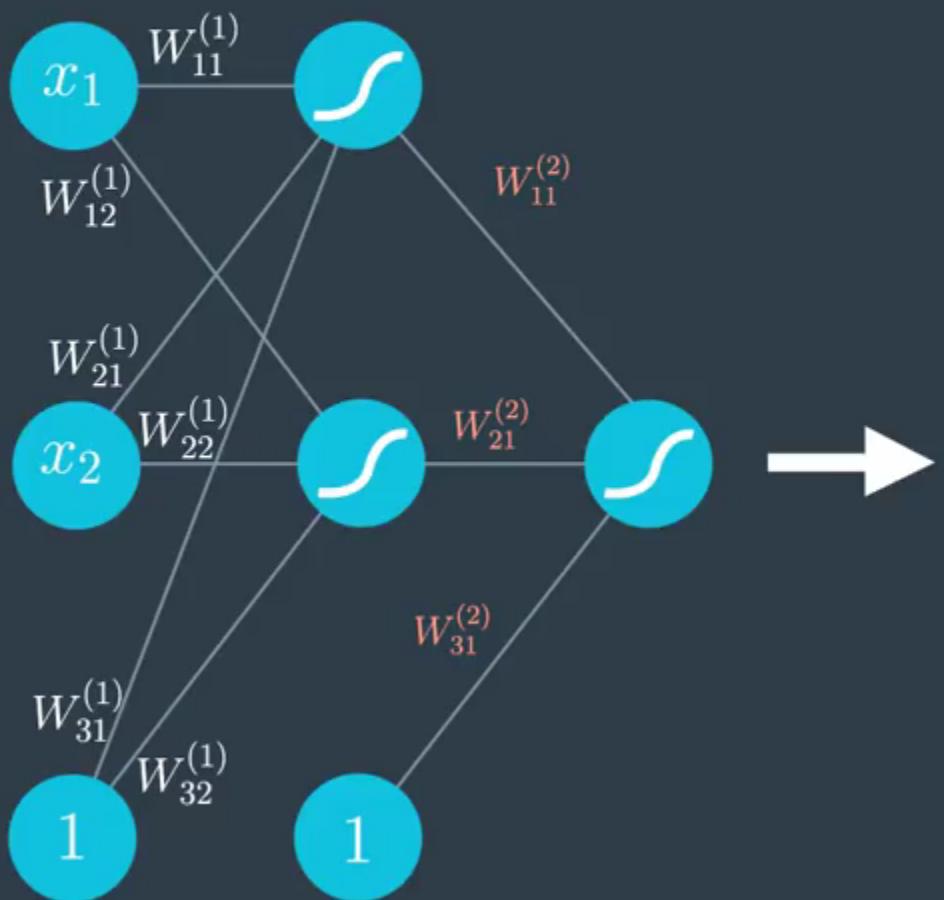
$$\begin{pmatrix} W_{11}^{(1)} & W_{12}^{(1)} \\ W_{21}^{(1)} & W_{22}^{(1)} \\ W_{31}^{(1)} & W_{32}^{(1)} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ 1 \end{pmatrix} \quad \begin{pmatrix} W_{11}^{(2)} & W_{12}^{(2)} \\ W_{21}^{(2)} & W_{31}^{(2)} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ 1 \end{pmatrix}$$

Feedforward



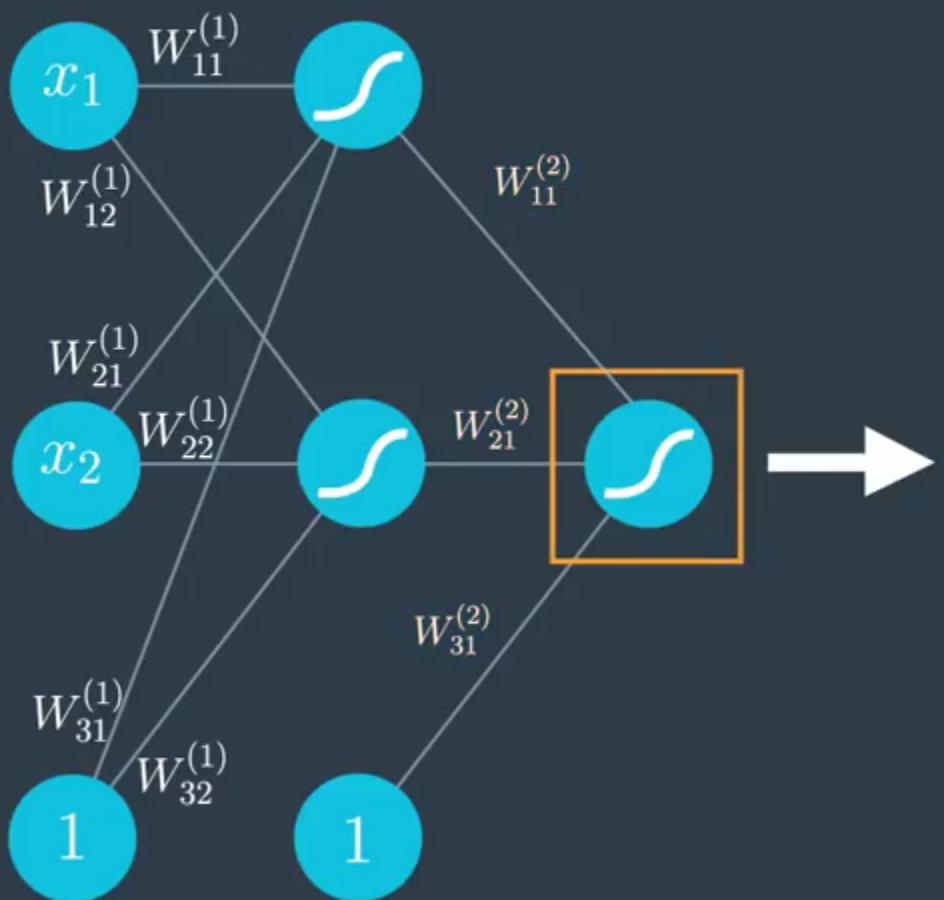
$$\sigma \begin{pmatrix} W_{11}^{(1)} & W_{12}^{(1)} \\ W_{21}^{(1)} & W_{22}^{(1)} \\ W_{31}^{(1)} & W_{32}^{(1)} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ 1 \end{pmatrix}$$

Feedforward



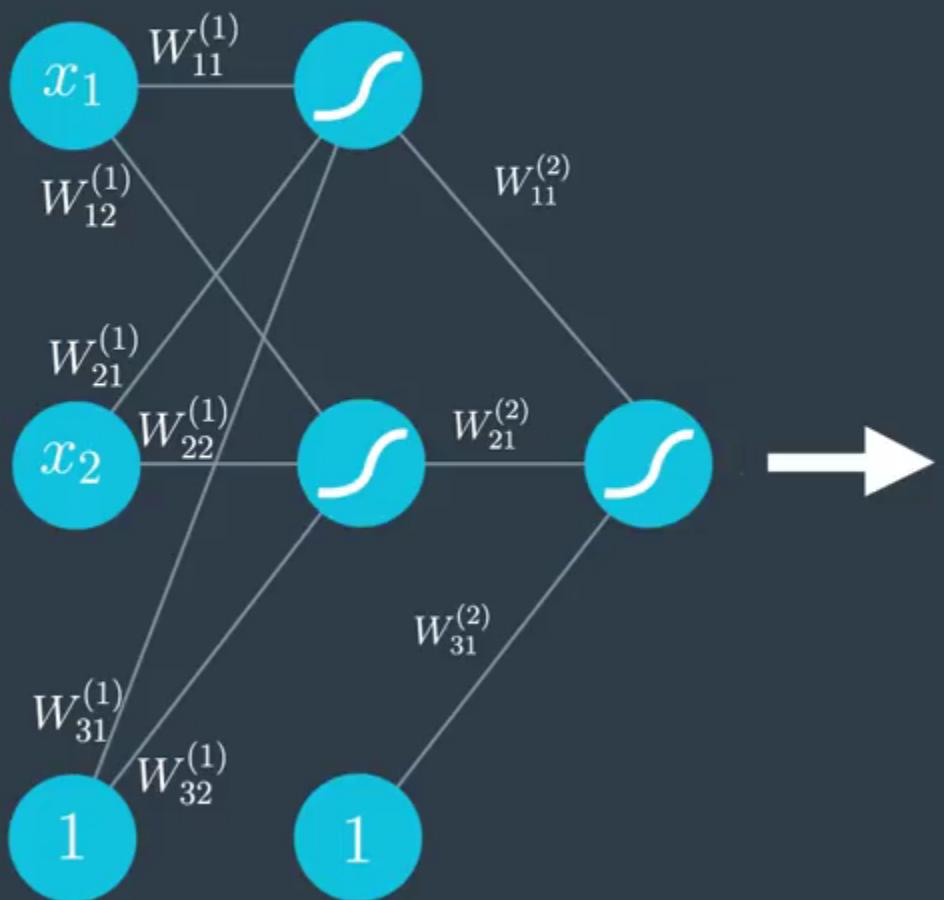
$$\begin{pmatrix} W_{11}^{(2)} \\ W_{21}^{(2)} \\ W_{31}^{(2)} \end{pmatrix} \sigma \begin{pmatrix} W_{11}^{(1)} & W_{12}^{(1)} \\ W_{21}^{(1)} & W_{22}^{(1)} \\ W_{31}^{(1)} & W_{32}^{(1)} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ 1 \end{pmatrix}$$

Feedforward



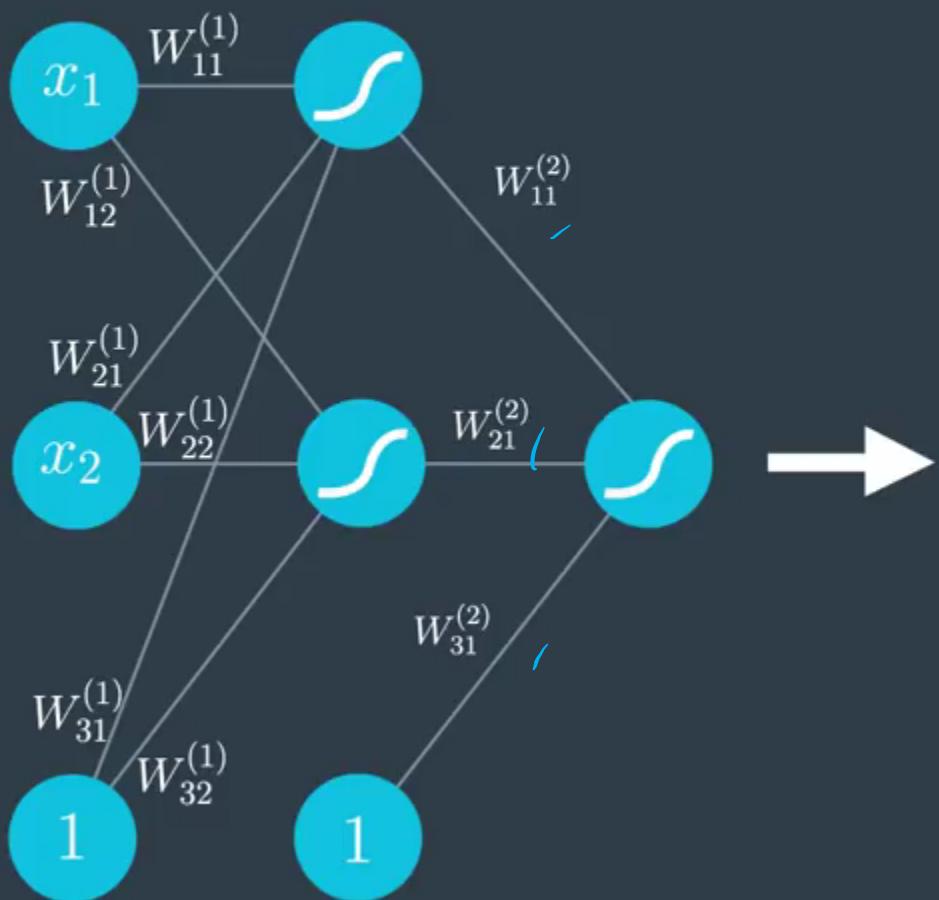
$$\sigma \begin{pmatrix} W_{11}^{(2)} \\ W_{21}^{(2)} \\ W_{31}^{(2)} \end{pmatrix} \sigma \begin{pmatrix} W_{11}^{(1)} & W_{12}^{(1)} \\ W_{21}^{(1)} & W_{22}^{(1)} \\ W_{31}^{(1)} & W_{32}^{(1)} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ 1 \end{pmatrix}$$

Feedforward



$$\hat{y} = \sigma \begin{pmatrix} W_{11}^{(2)} \\ W_{21}^{(2)} \\ W_{31}^{(2)} \end{pmatrix} \sigma \begin{pmatrix} W_{11}^{(1)} & W_{12}^{(1)} \\ W_{21}^{(1)} & W_{22}^{(1)} \\ W_{31}^{(1)} & W_{32}^{(1)} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ 1 \end{pmatrix}$$

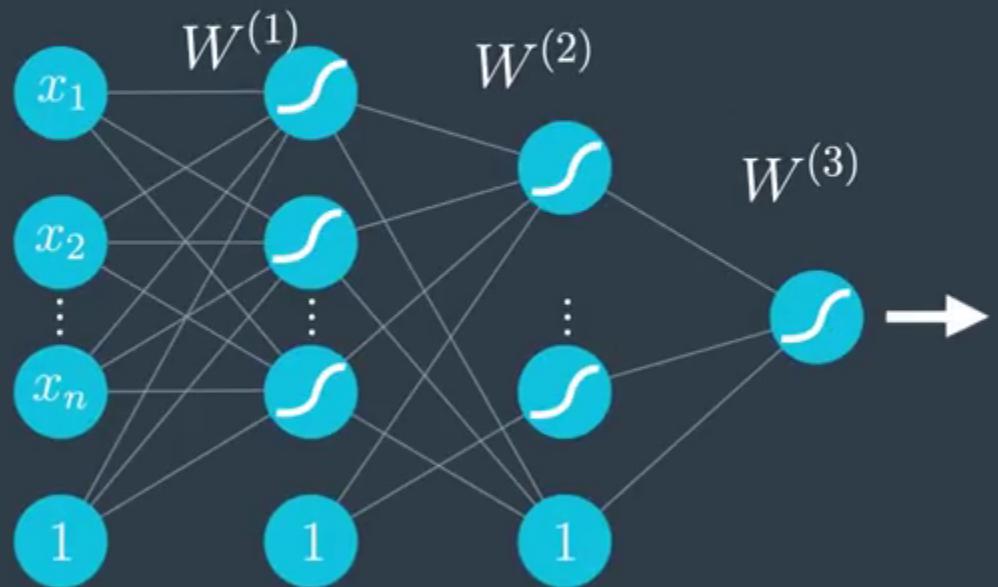
Feedforward



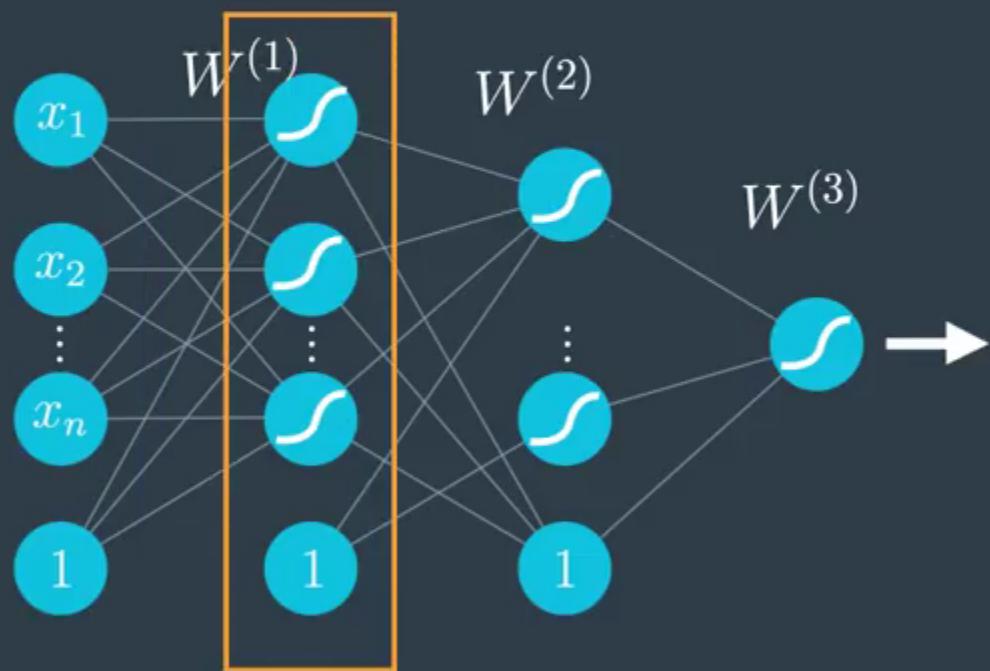
$$\hat{y} = \sigma \begin{pmatrix} W_{11}^{(2)} \\ W_{21}^{(2)} \\ W_{31}^{(2)} \end{pmatrix} \sigma \begin{pmatrix} W_{11}^{(1)} & W_{12}^{(1)} \\ W_{21}^{(1)} & W_{22}^{(1)} \\ W_{31}^{(1)} & W_{32}^{(1)} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ 1 \end{pmatrix}$$

$$\hat{y} = \sigma \circ W^{(2)} \circ \sigma \circ W^{(1)}(x)$$

Multi-layer Perceptron



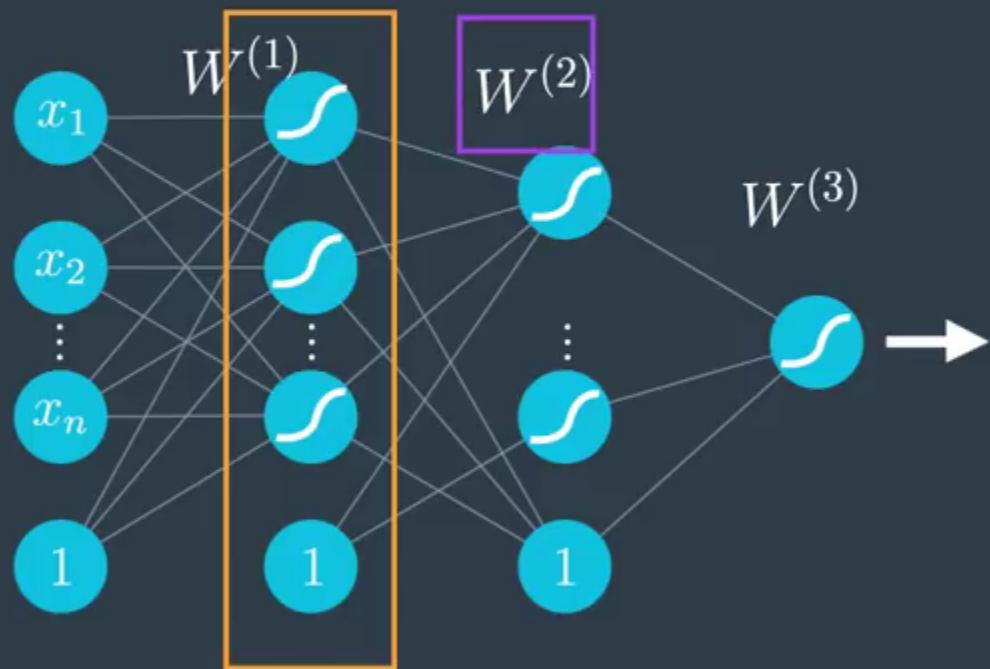
Multi-layer Perceptron



PREDICTION

$$\hat{y} = \sigma \circ W^{(1)}(x)$$

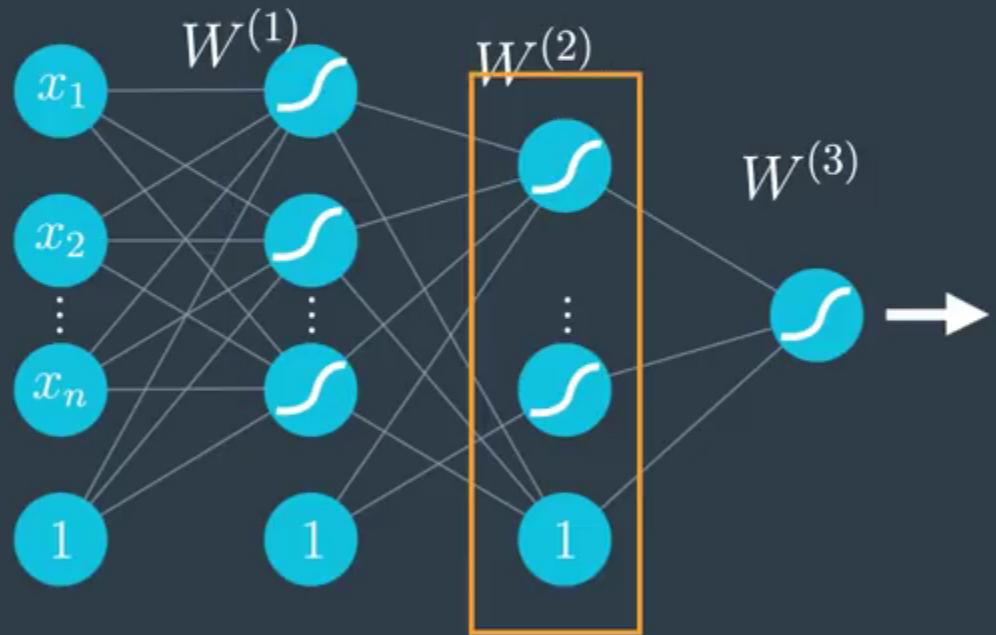
Multi-layer Perceptron



PREDICTION

$$\hat{y} = W^{(2)} \circ \sigma \circ W^{(1)}(x)$$

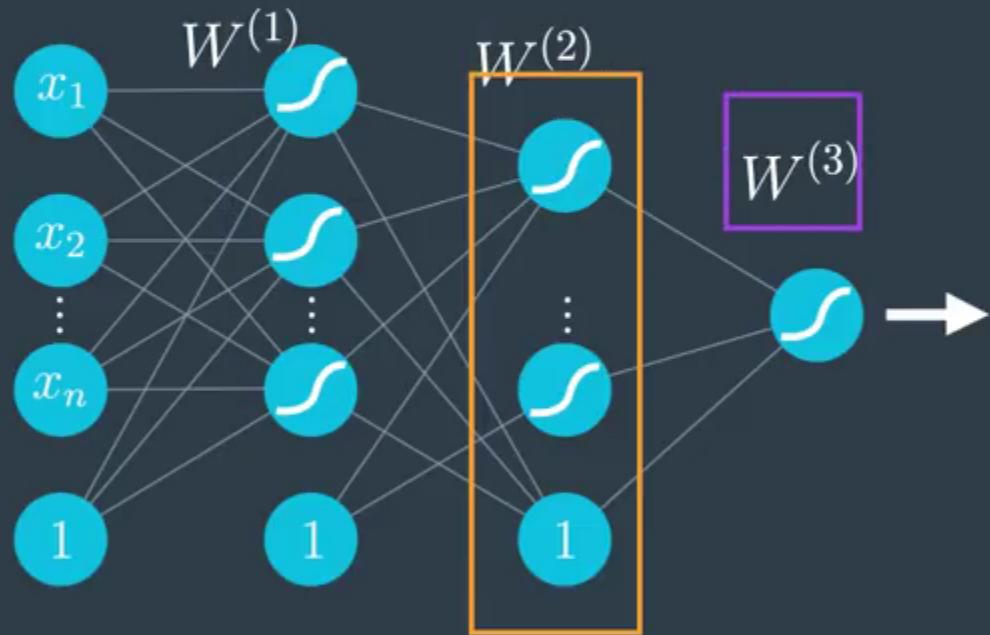
Multi-layer Perceptron



PREDICTION

$$\hat{y} = \sigma \circ W^{(2)} \circ \sigma \circ W^{(1)}(x)$$

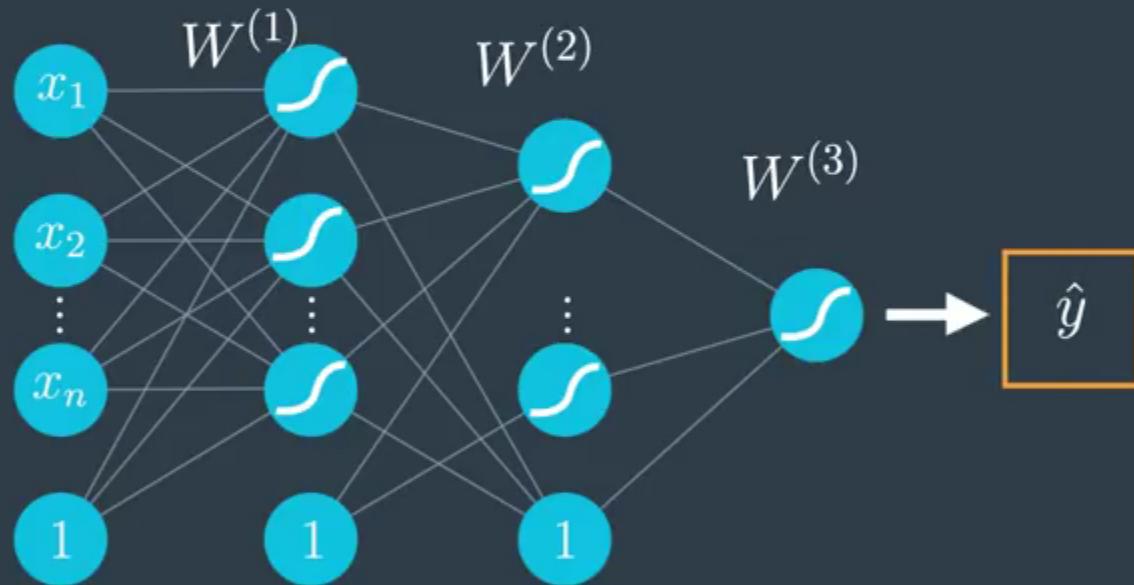
Multi-layer Perceptron



PREDICTION

$$\hat{y} = \sigma \circ W^{(3)} \circ \sigma \circ W^{(2)} \circ \sigma \circ W^{(1)}(x)$$

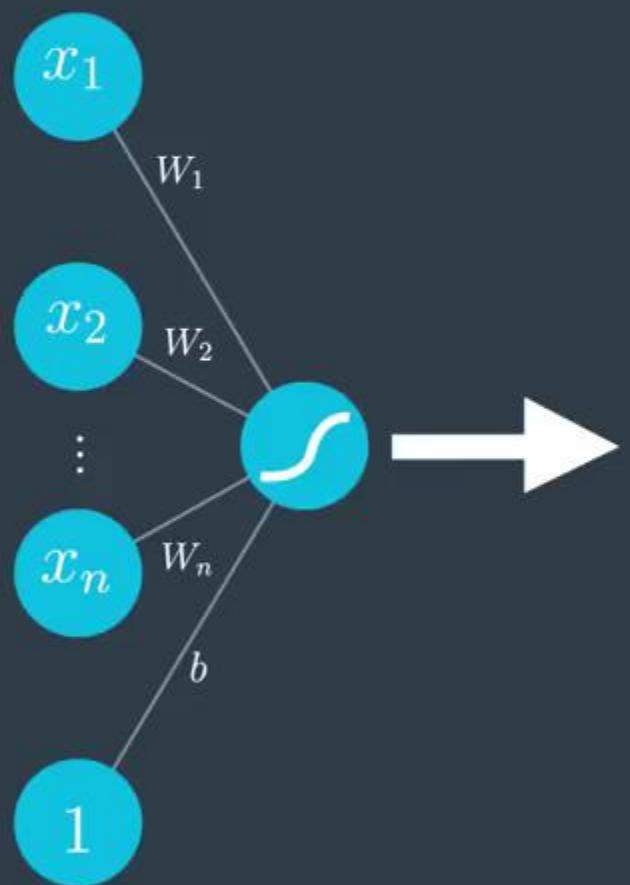
Multi-layer Perceptron



PREDICTION

$$\hat{y} = \sigma \circ W^{(3)} \circ \sigma \circ W^{(2)} \circ \sigma \circ W^{(1)}(x)$$

Perceptron

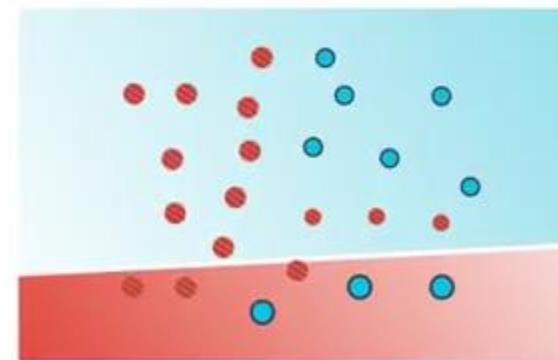


PREDICTION

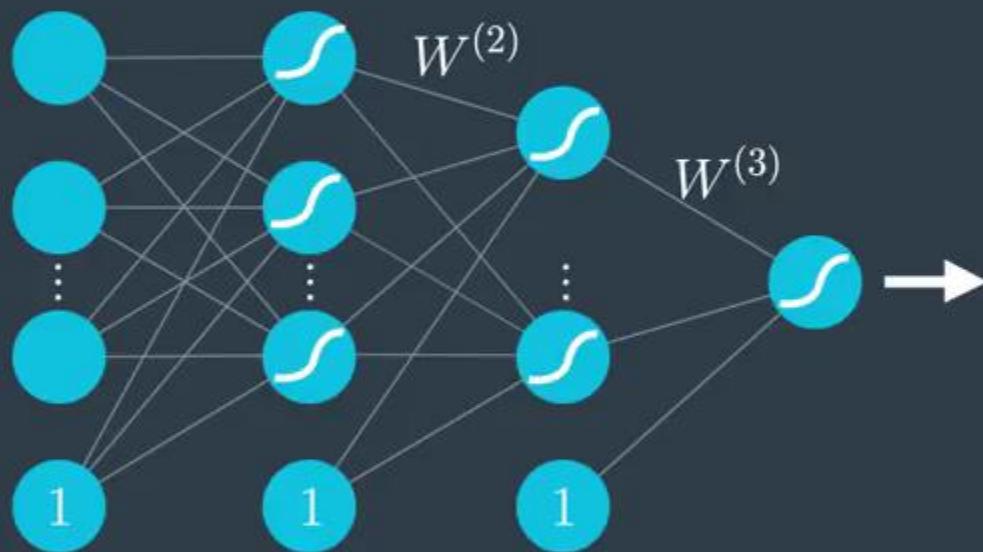
$$\hat{y} = \sigma(Wx + b)$$

ERROR FUNCTION

$$E(W) = -\frac{1}{m} \sum_{i=1}^m y_i \ln(\hat{y}_i) + (1 - y_i) \ln(1 - \hat{y}_i)$$



Multi-layer Perceptron

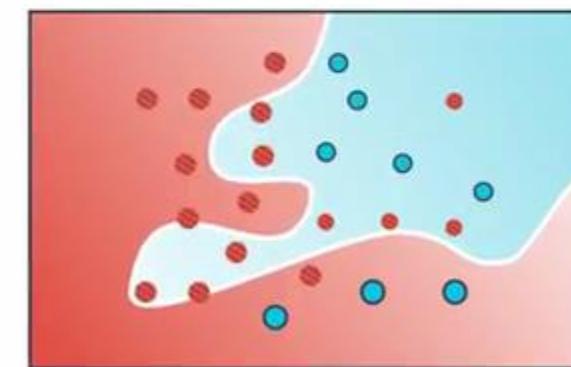


PREDICTION

$$\hat{y} = \sigma \circ W^{(3)} \circ \sigma \circ W^{(2)} \circ \sigma \circ W^{(1)}(x)$$

ERROR FUNCTION

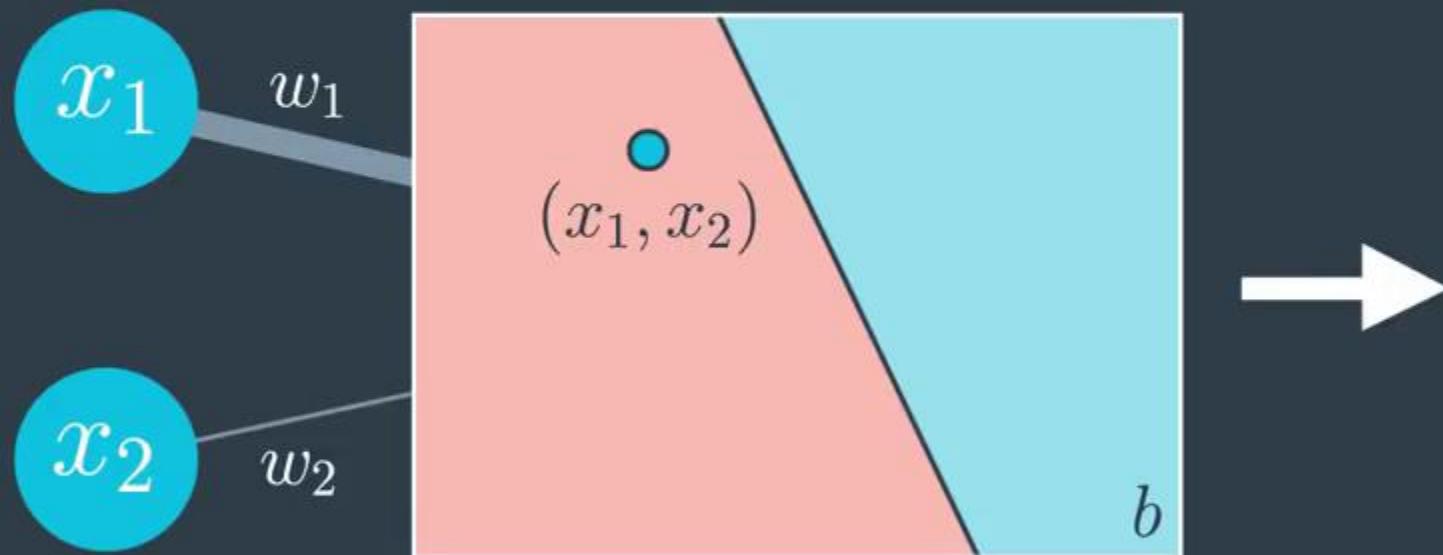
$$E(W) = -\frac{1}{m} \sum_{i=1}^m y_i \ln(\hat{y}_i) + (1 - y_i) \ln(1 - \hat{y}_i)$$



Backpropagation

$$x = (x_1, x_2)$$

$$y = 1$$

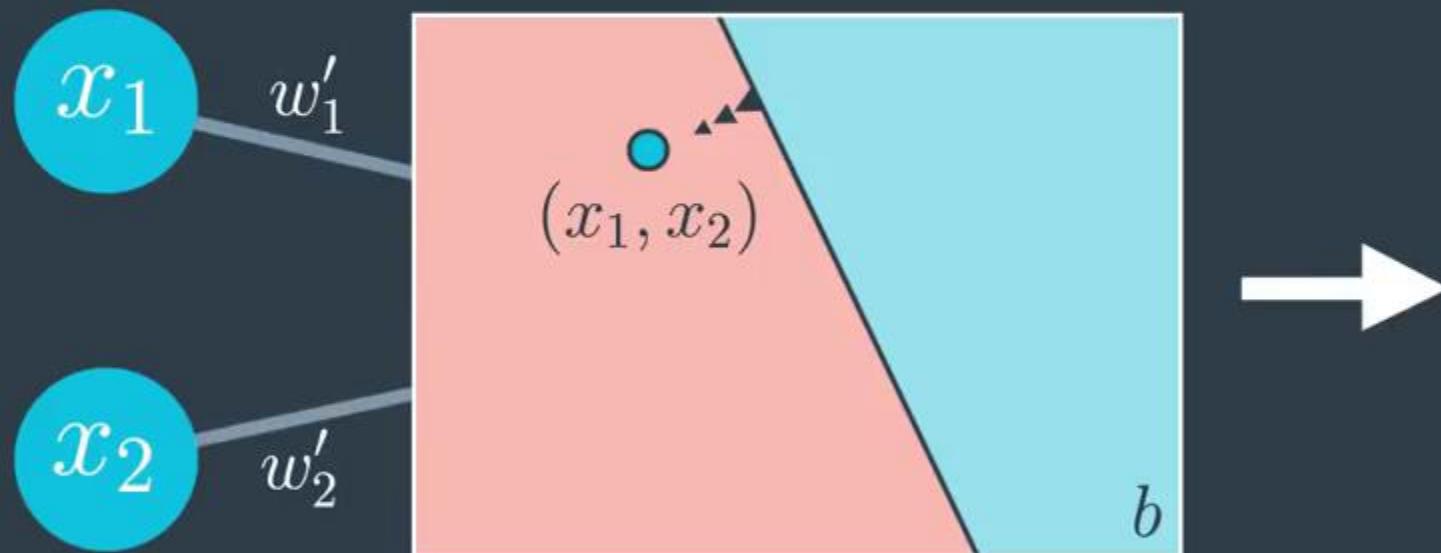


$$w_1x_1 + w_2x_2 + b$$

Backpropagation

$$x = (x_1, x_2)$$

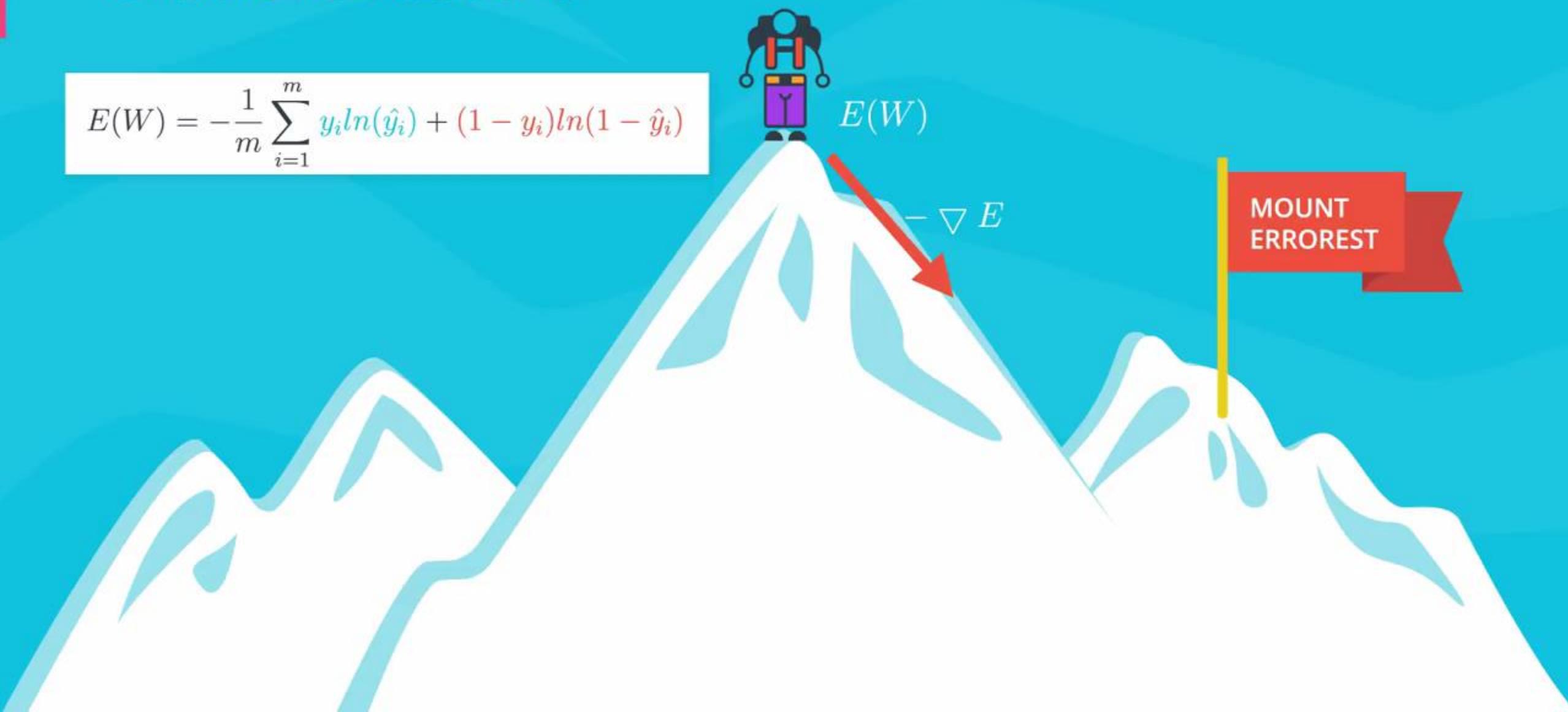
$$y = 1$$



$$w_1x_1 + w_2x_2 + b$$

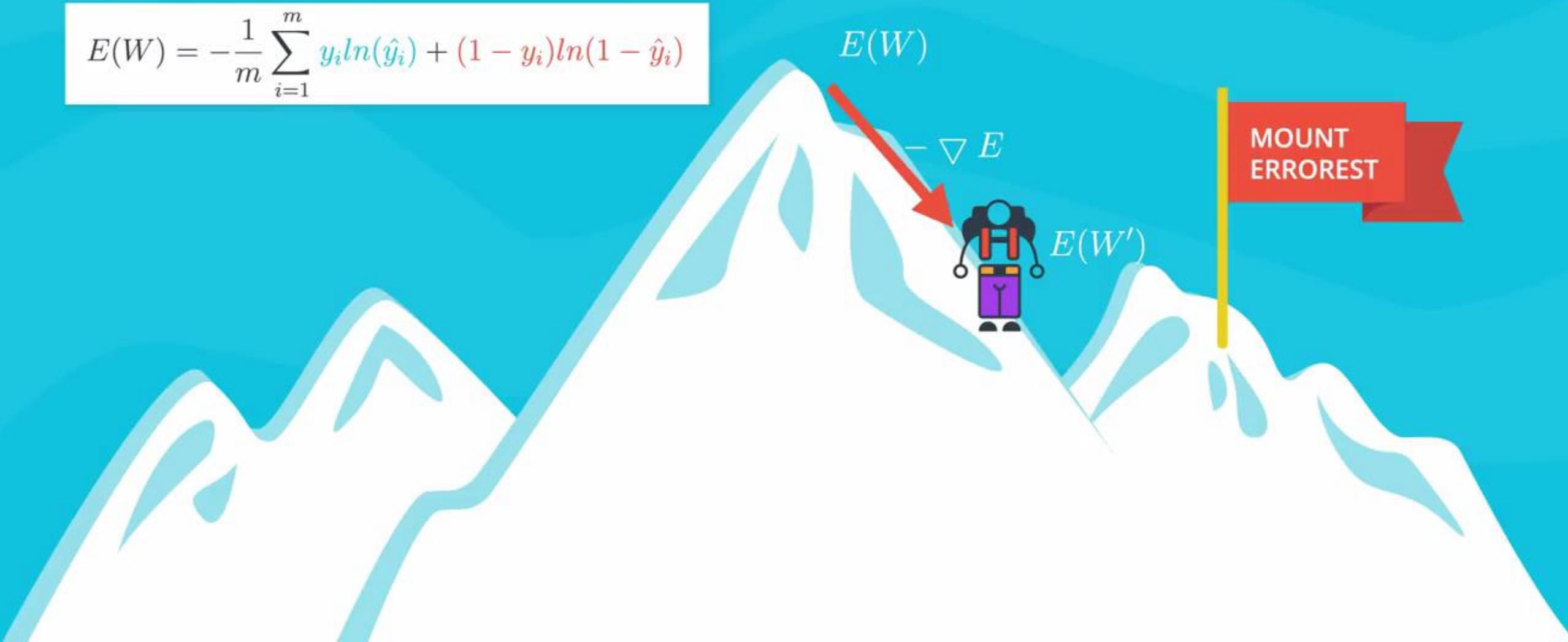
Gradient Descent

$$E(W) = -\frac{1}{m} \sum_{i=1}^m y_i \ln(\hat{y}_i) + (1 - y_i) \ln(1 - \hat{y}_i)$$



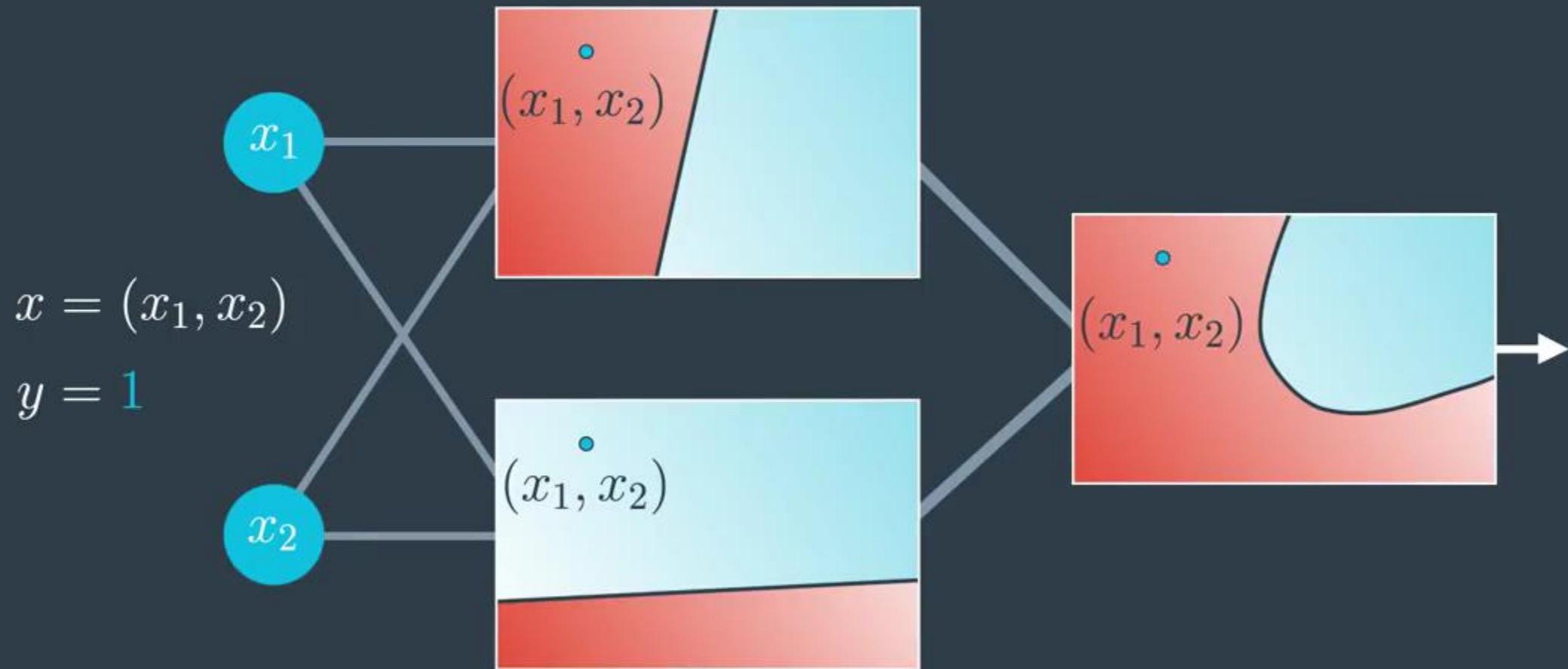
Gradient Descent

$$E(W) = -\frac{1}{m} \sum_{i=1}^m y_i \ln(\hat{y}_i) + (1 - y_i) \ln(1 - \hat{y}_i)$$

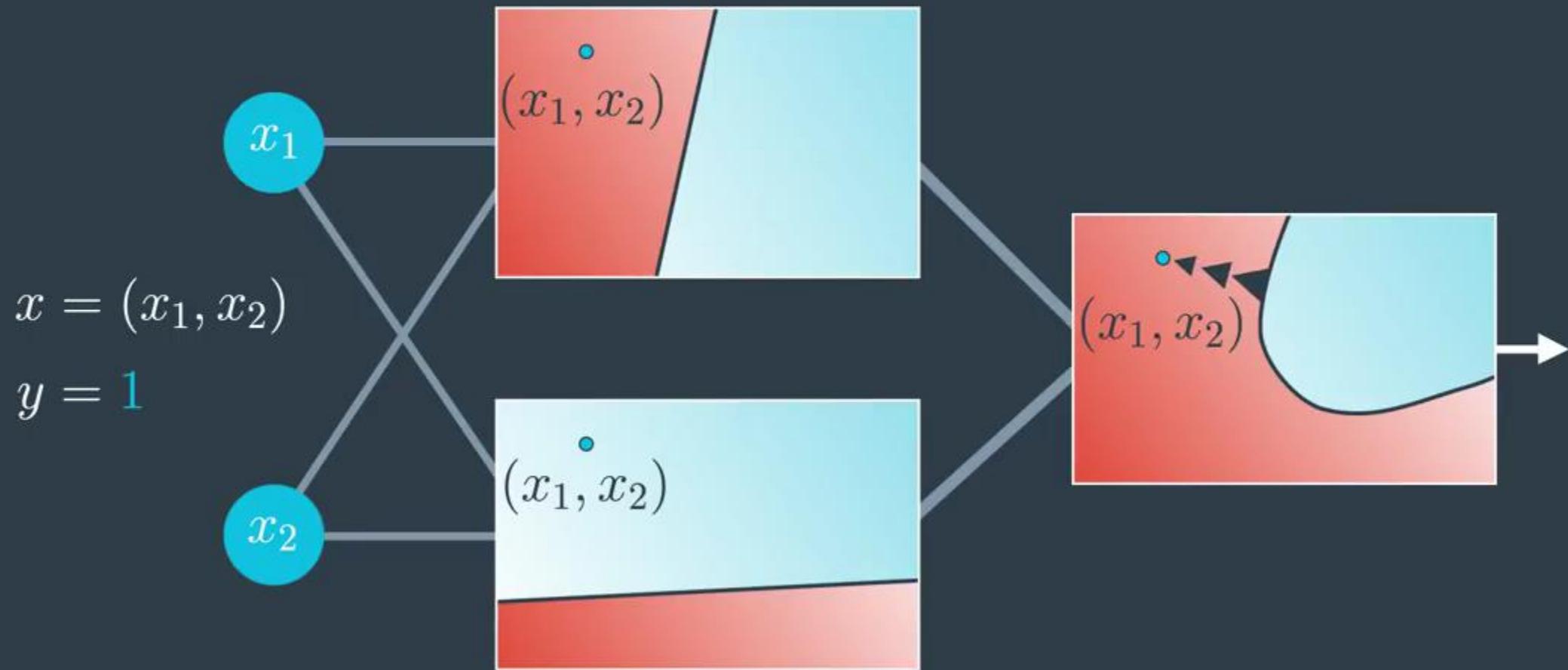


What about
Multi-Layer Perceptrons?

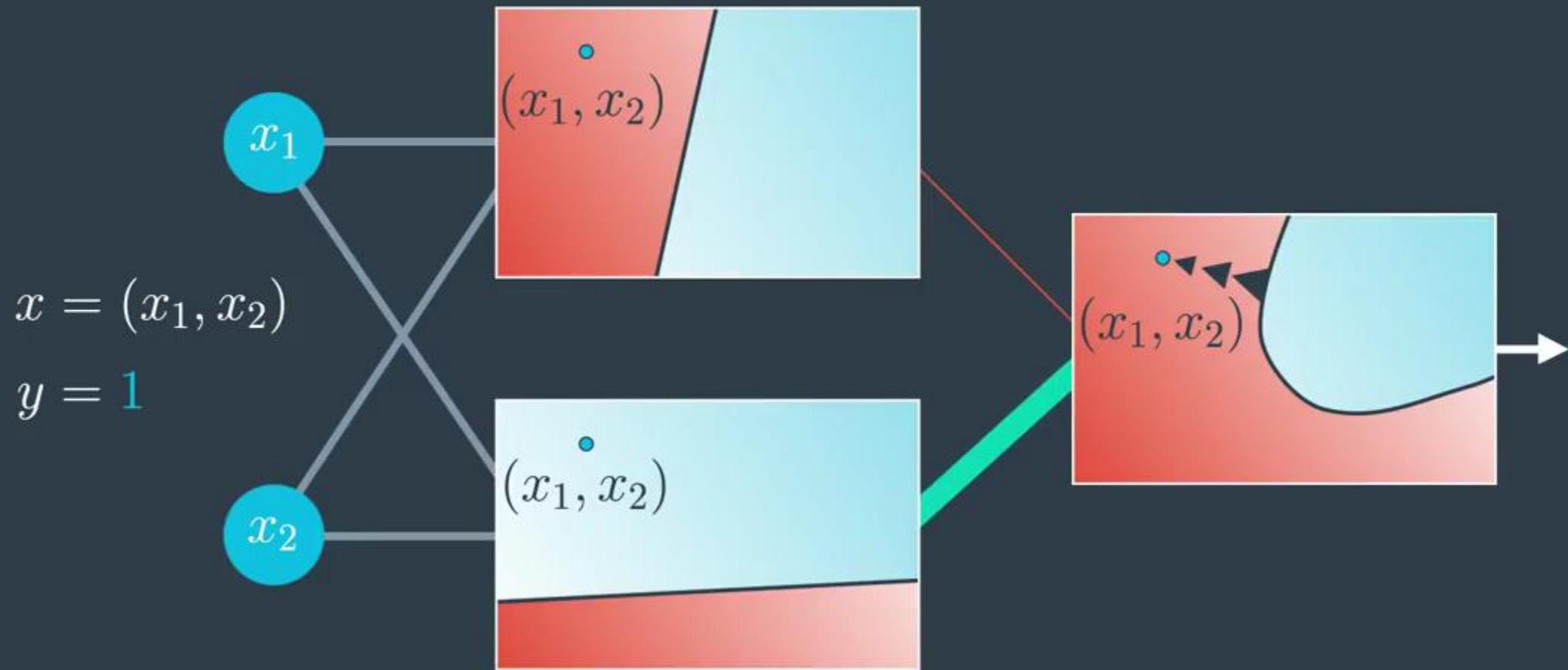
FeedForward



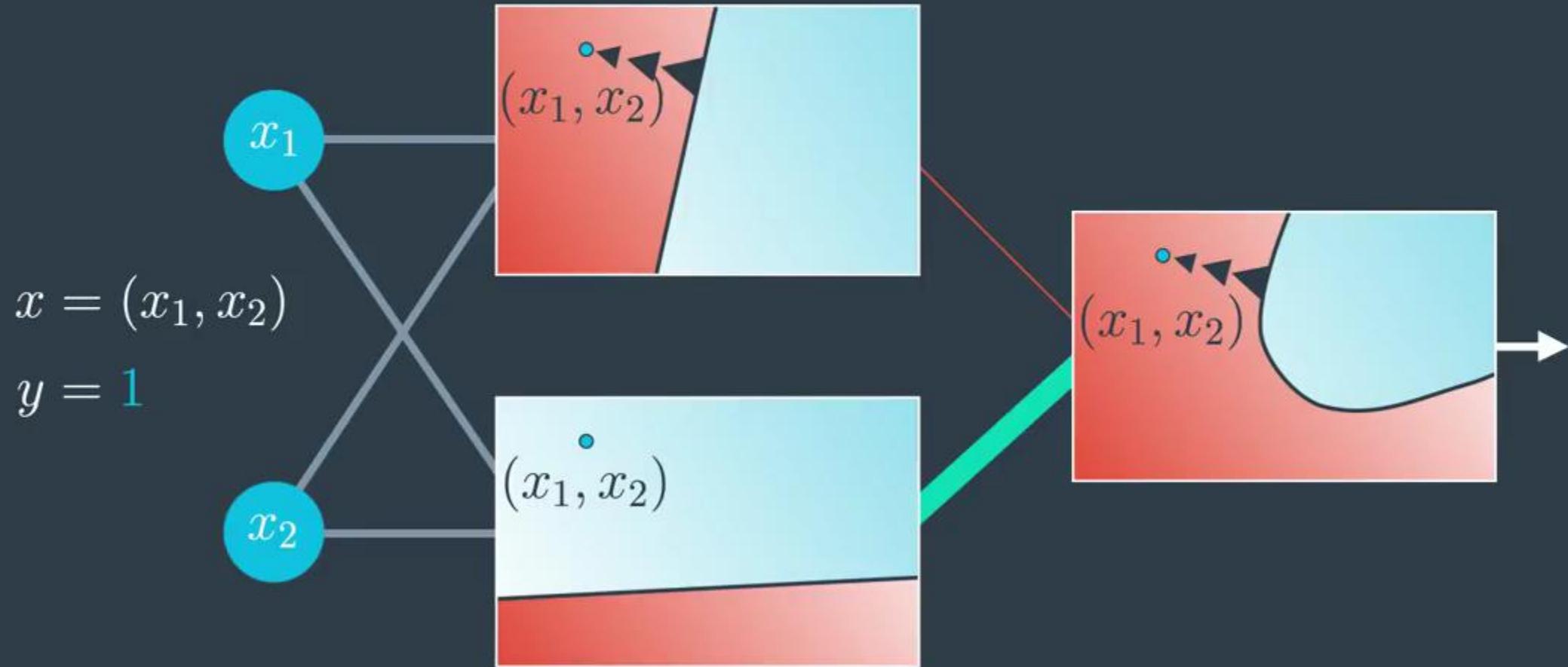
Backpropagation



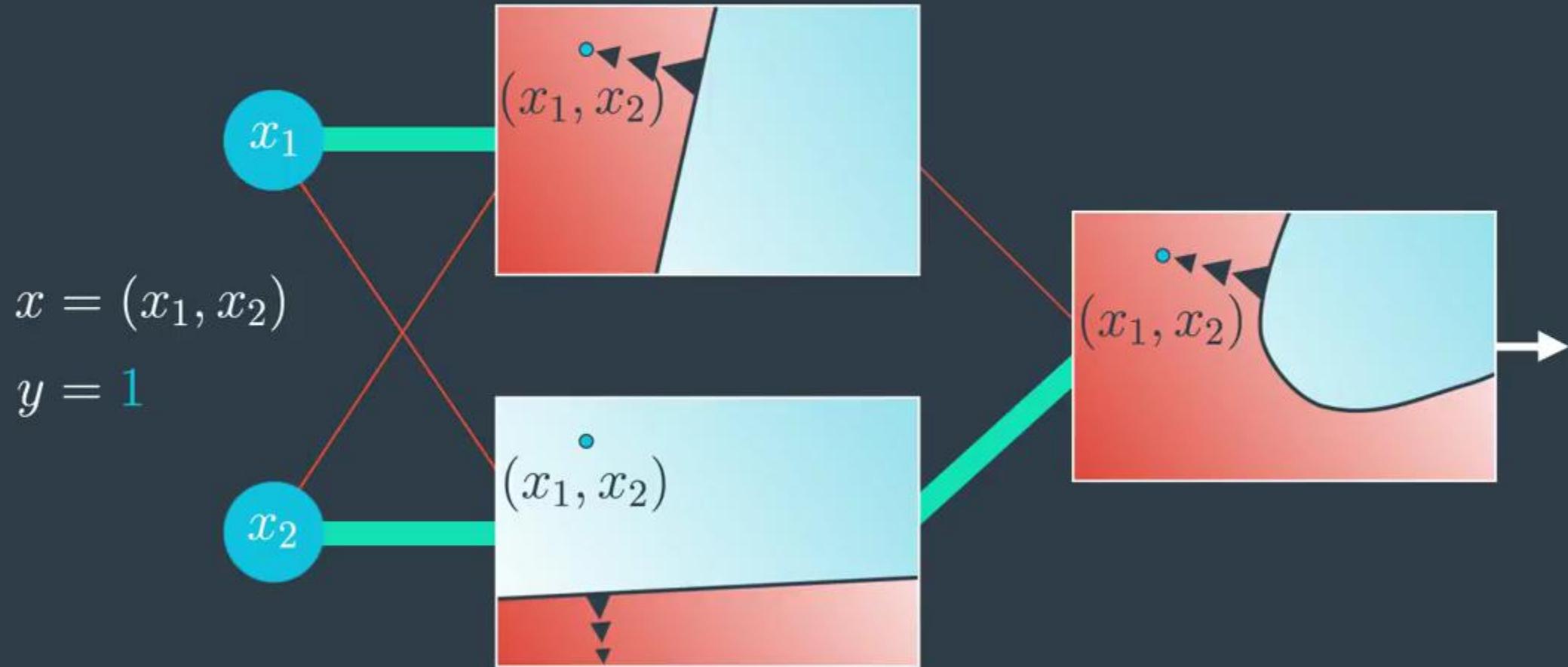
Backpropagation



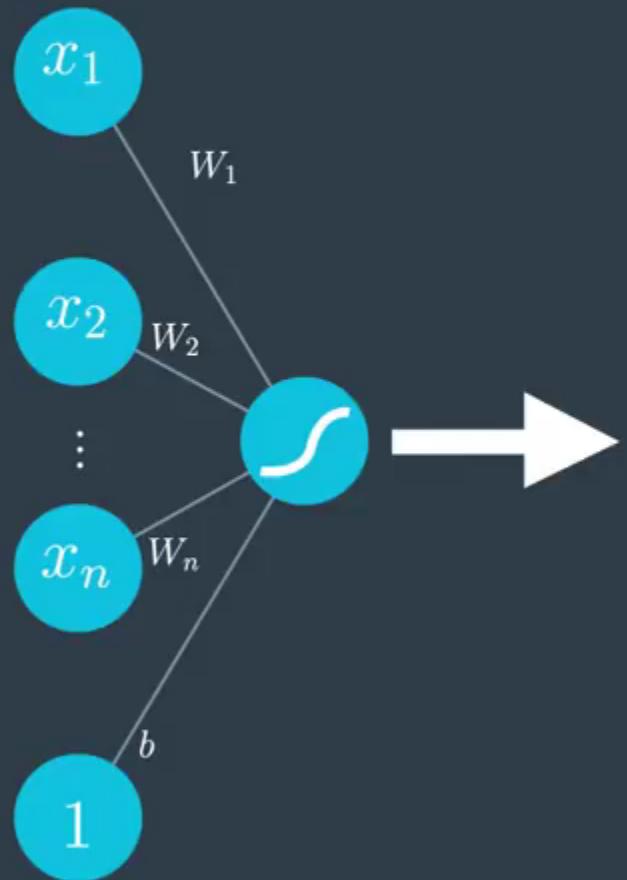
Backpropagation



Backpropagation



Perceptron



PREDICTION

$$\hat{y} = \sigma(Wx + b)$$

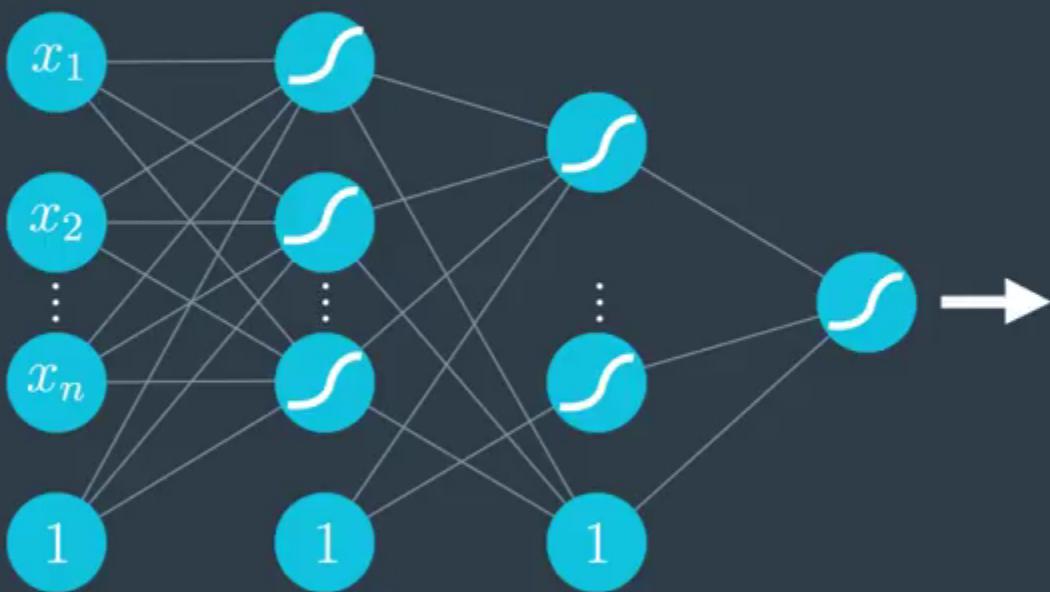
ERROR FUNCTION

$$E(W) = -\frac{1}{m} \sum_{i=1}^m y_i \ln(\hat{y}_i) + (1 - y_i) \ln(1 - \hat{y}_i)$$

GRADIENT OF THE ERROR FUNCTION

$$\nabla E = \left(\frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n}, \frac{\partial E}{\partial b} \right)$$

Multi-layer Perceptron



PREDICTION

$$\hat{y} = \sigma W^{(3)} \circ \sigma W^{(2)} \circ \sigma \circ W^{(1)}(x)$$

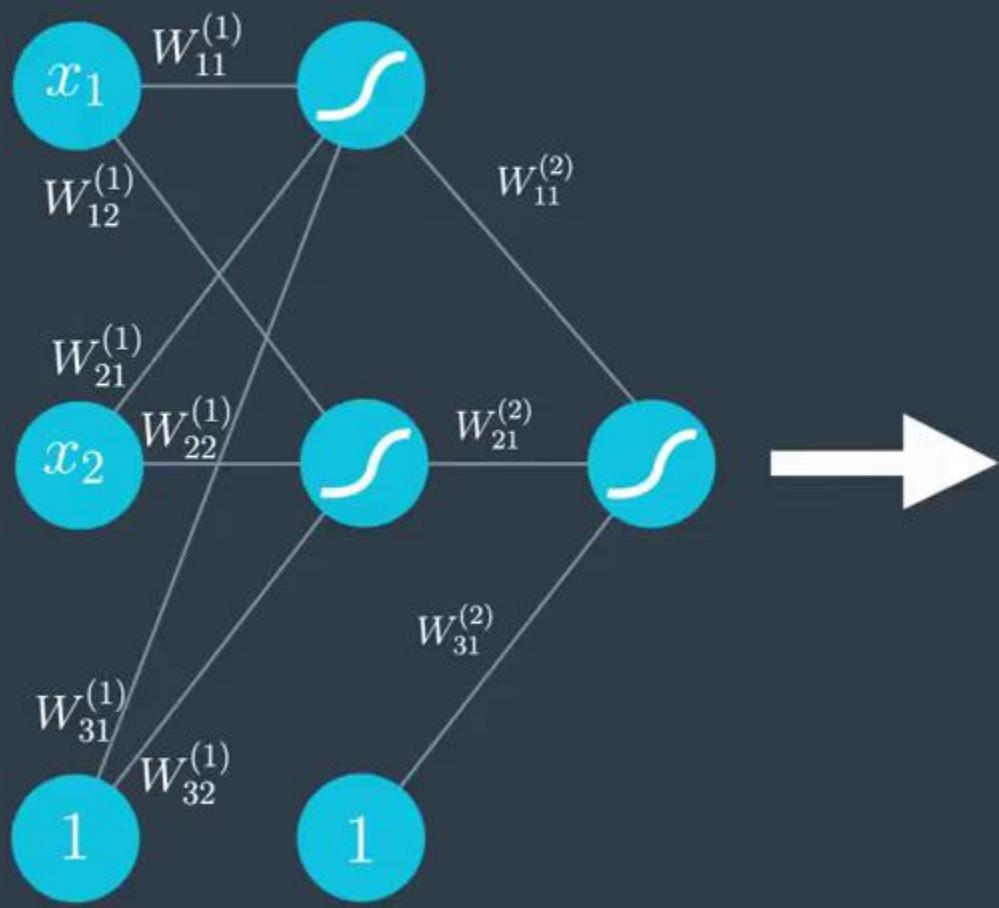
ERROR FUNCTION

$$E(W) = -\frac{1}{m} \sum_{i=1}^m y_i \ln(\hat{y}_i) + (1 - y_i) \ln(1 - \hat{y}_i)$$

GRADIENT OF THE ERROR FUNCTION

$$\nabla E = \left(\dots, \frac{\partial E}{\partial w_j^{(i)}}, \dots \right)$$

Backpropagation



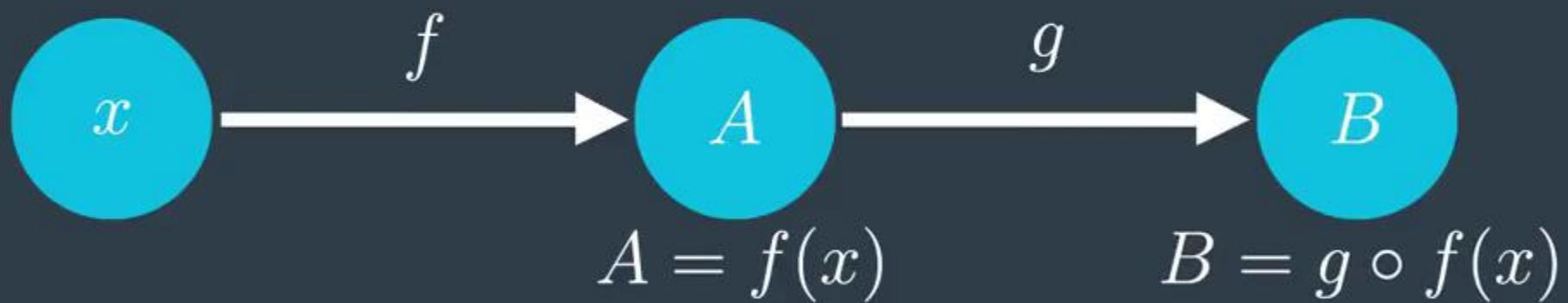
$$\hat{y} = \sigma W^{(2)} \circ \sigma \circ W^{(1)}(x)$$

$$W^{(1)} = \begin{pmatrix} W_{11}^{(1)} & W_{12}^{(1)} \\ W_{21}^{(1)} & W_{22}^{(1)} \\ W_{31}^{(1)} & W_{32}^{(1)} \end{pmatrix} \quad W^{(2)} = \begin{pmatrix} W_{11}^{(2)} \\ W_{21}^{(2)} \\ W_{31}^{(2)} \end{pmatrix}$$

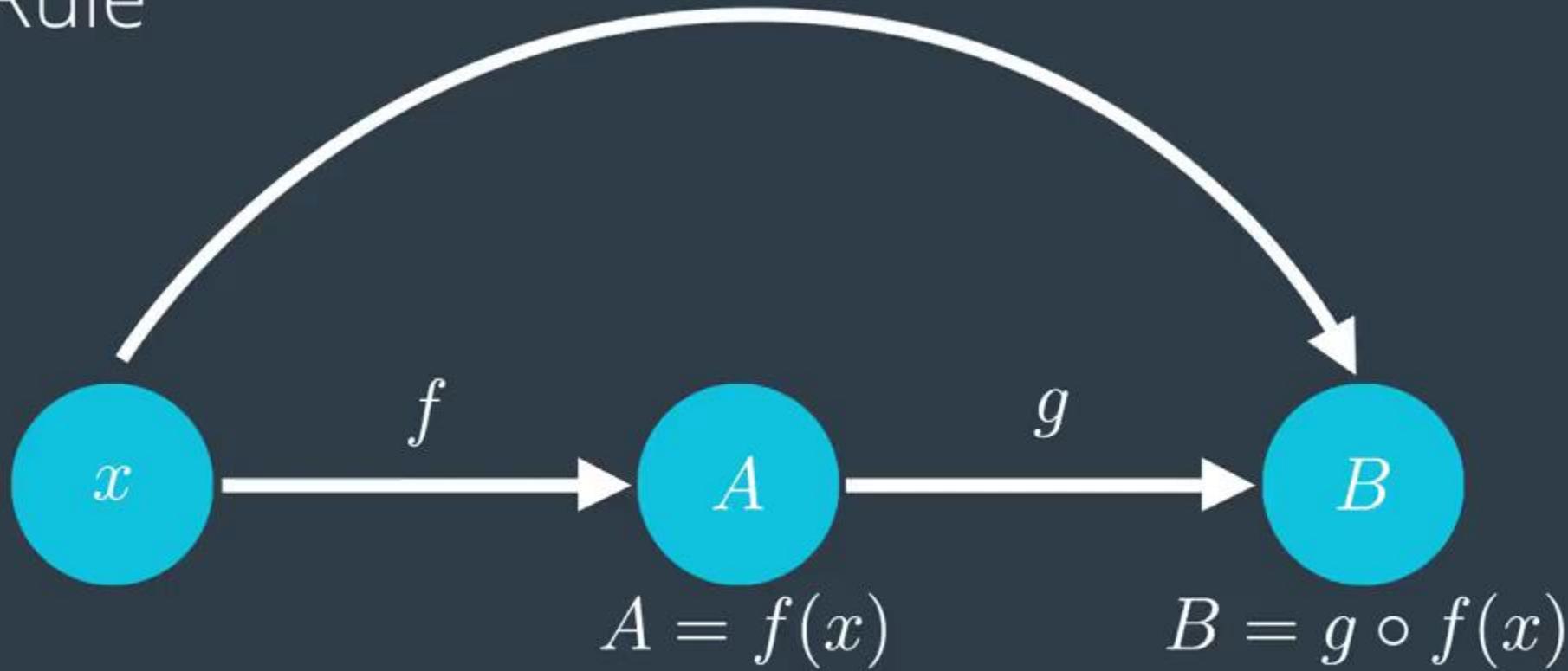
$$\nabla E = \begin{pmatrix} \frac{\partial E}{\partial W_{11}^{(1)}} & \frac{\partial E}{\partial W_{12}^{(1)}} & \frac{\partial E}{\partial W_{11}^{(2)}} \\ \frac{\partial E}{\partial W_{21}^{(1)}} & \frac{\partial E}{\partial W_{22}^{(1)}} & \frac{\partial E}{\partial W_{21}^{(2)}} \\ \frac{\partial E}{\partial W_{31}^{(1)}} & \frac{\partial E}{\partial W_{32}^{(1)}} & \frac{\partial E}{\partial W_{31}^{(2)}} \end{pmatrix}$$

$$W'_{ij}^{(k)} \leftarrow W_{ij}^{(k)} - \alpha \frac{\partial E}{\partial W_{ij}^{(k)}}$$

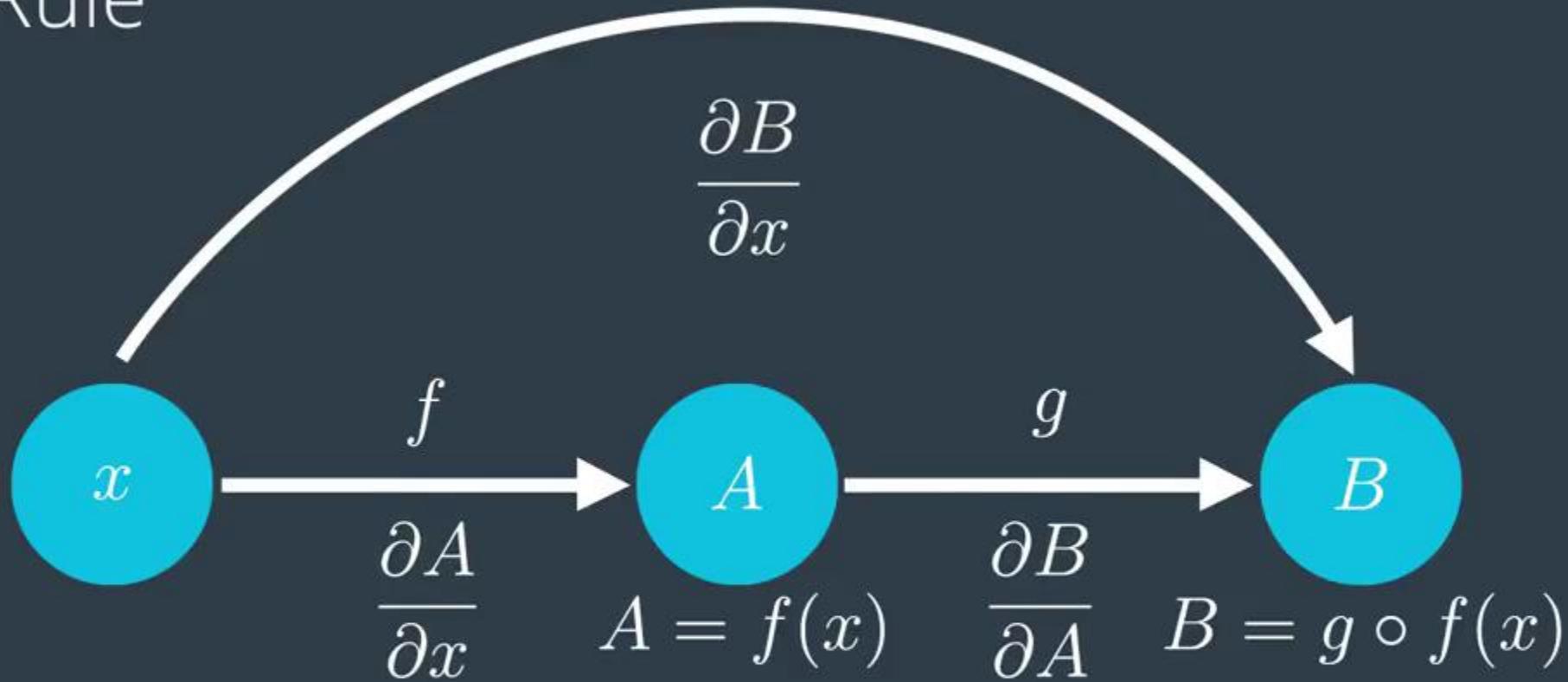
Chain Rule



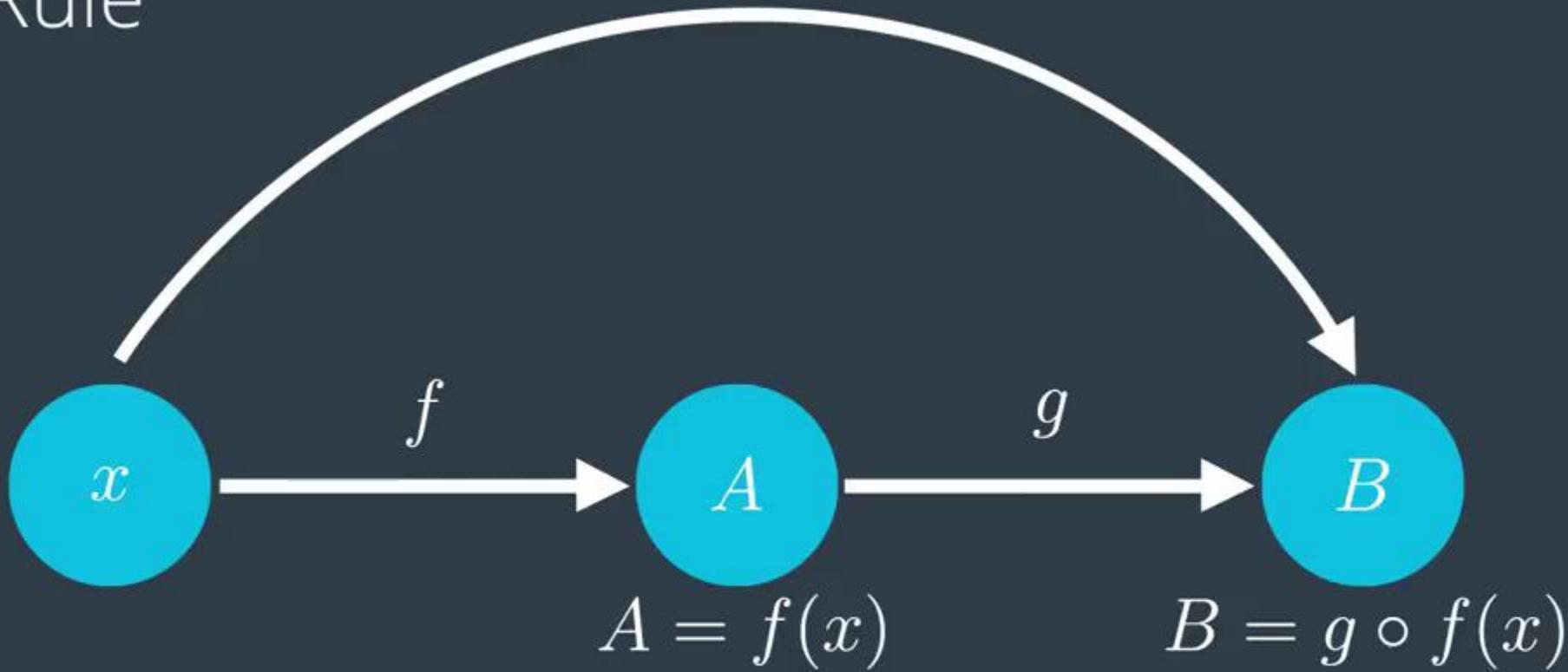
Chain Rule



Chain Rule

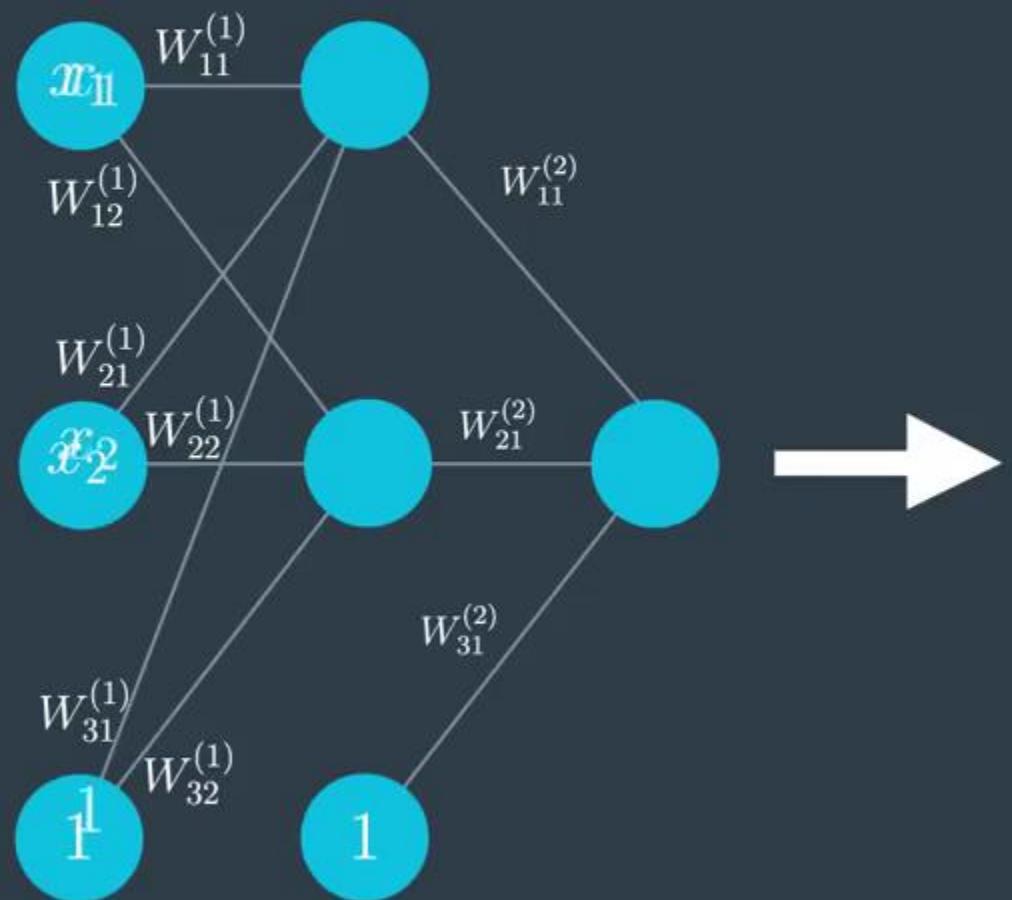


Chain Rule



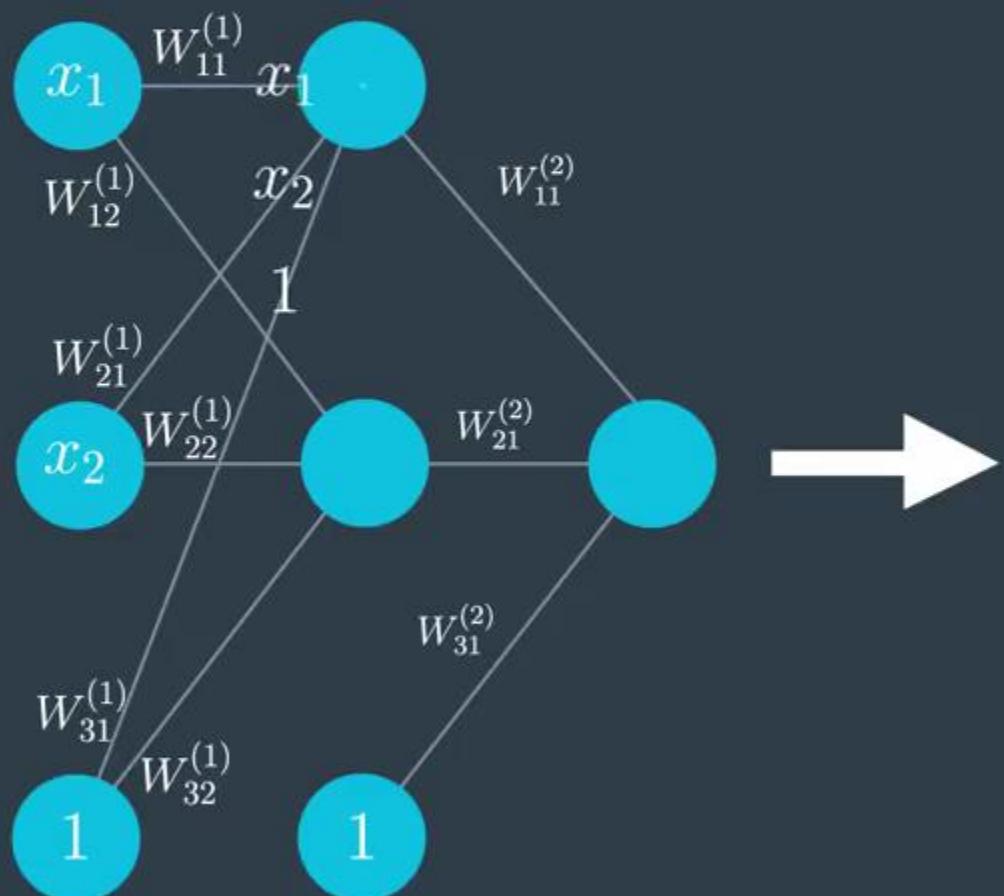
$$\frac{\partial B}{\partial x} = \frac{\partial B}{\partial A} \frac{\partial A}{\partial x}$$

Feedforward



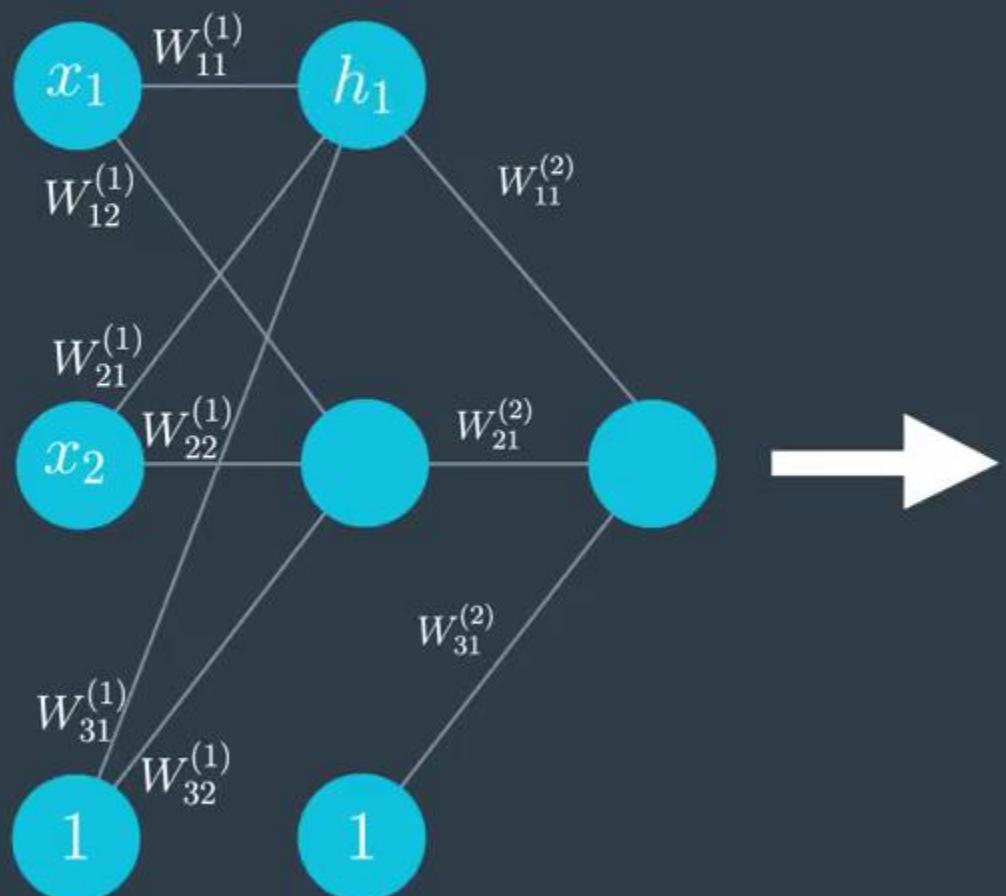
$$h_1 = W_{11}^{(1)}x_1 + W_{21}^{(1)}x_2 + W_{31}^{(1)}$$

Feedforward



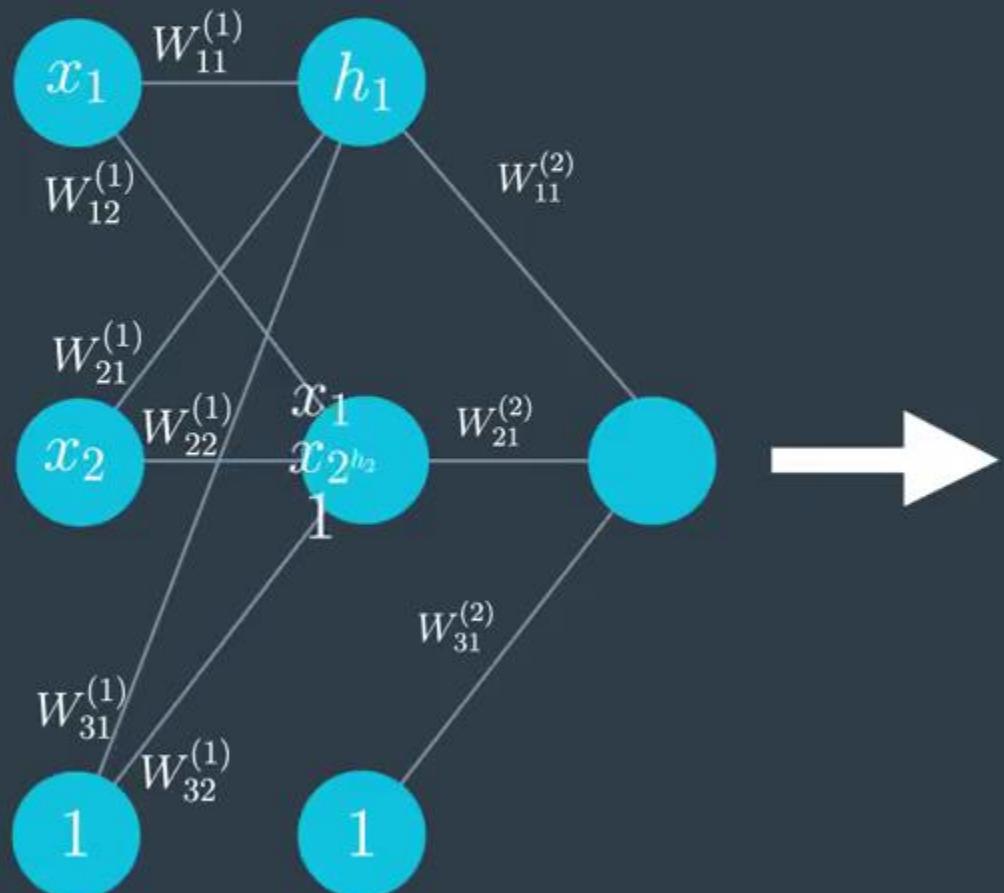
$$h_1 = W_{11}^{(1)}x_1 + W_{21}^{(1)}x_2 + W_{31}^{(1)}$$

Feedforward



$$h_1 = W_{11}^{(1)}x_1 + W_{21}^{(1)}x_2 + W_{31}^{(1)}$$

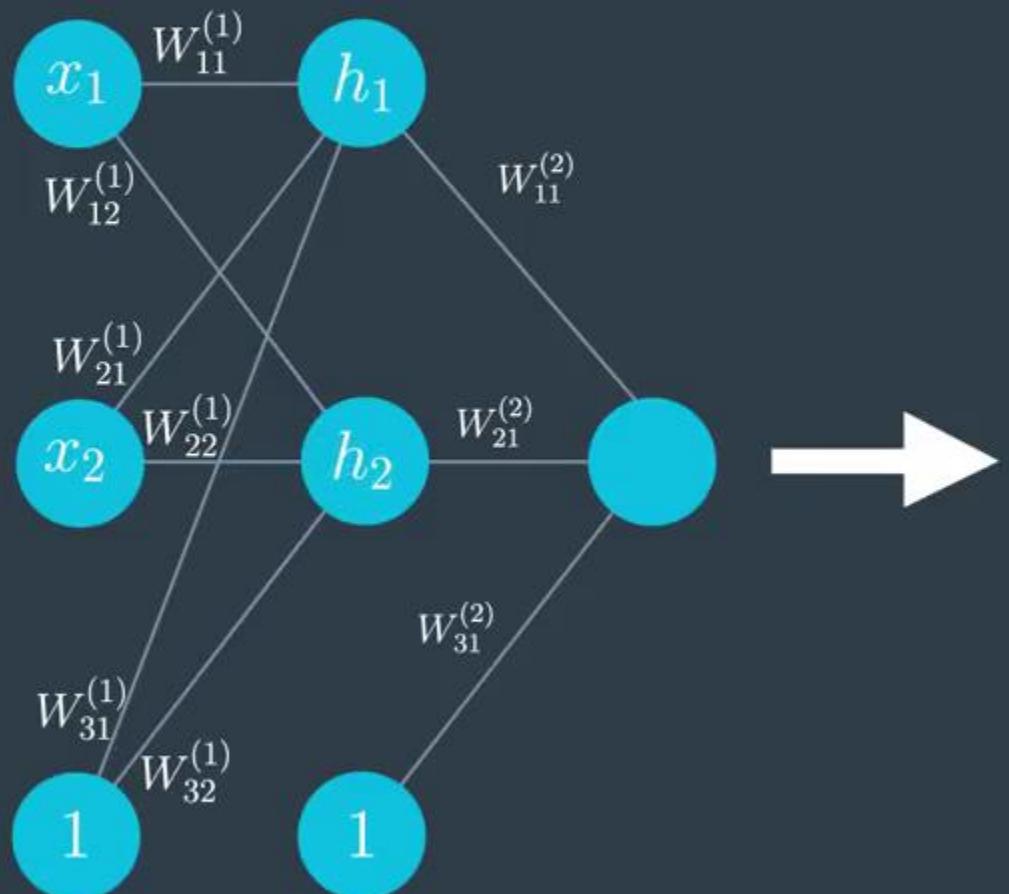
Feedforward



$$h_1 = W_{11}^{(1)}x_1 + W_{21}^{(1)}x_2 + W_{31}^{(1)}$$

$$h_2 = W_{12}^{(1)}x_1 + W_{22}^{(1)}x_2 + W_{32}^{(1)}$$

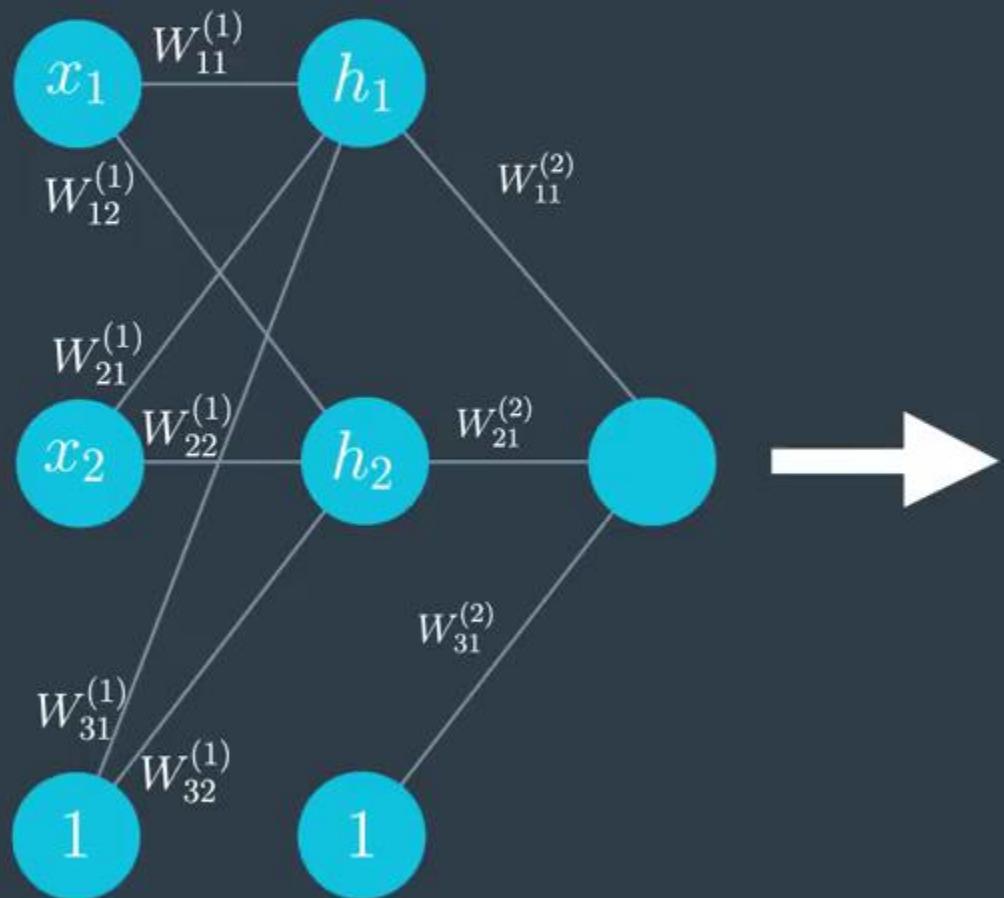
Feedforward



$$h_1 = W_{11}^{(1)}x_1 + W_{21}^{(1)}x_2 + W_{31}^{(1)}$$

$$h_2 = W_{12}^{(1)}x_1 + W_{22}^{(1)}x_2 + W_{32}^{(1)}$$

Feedforward

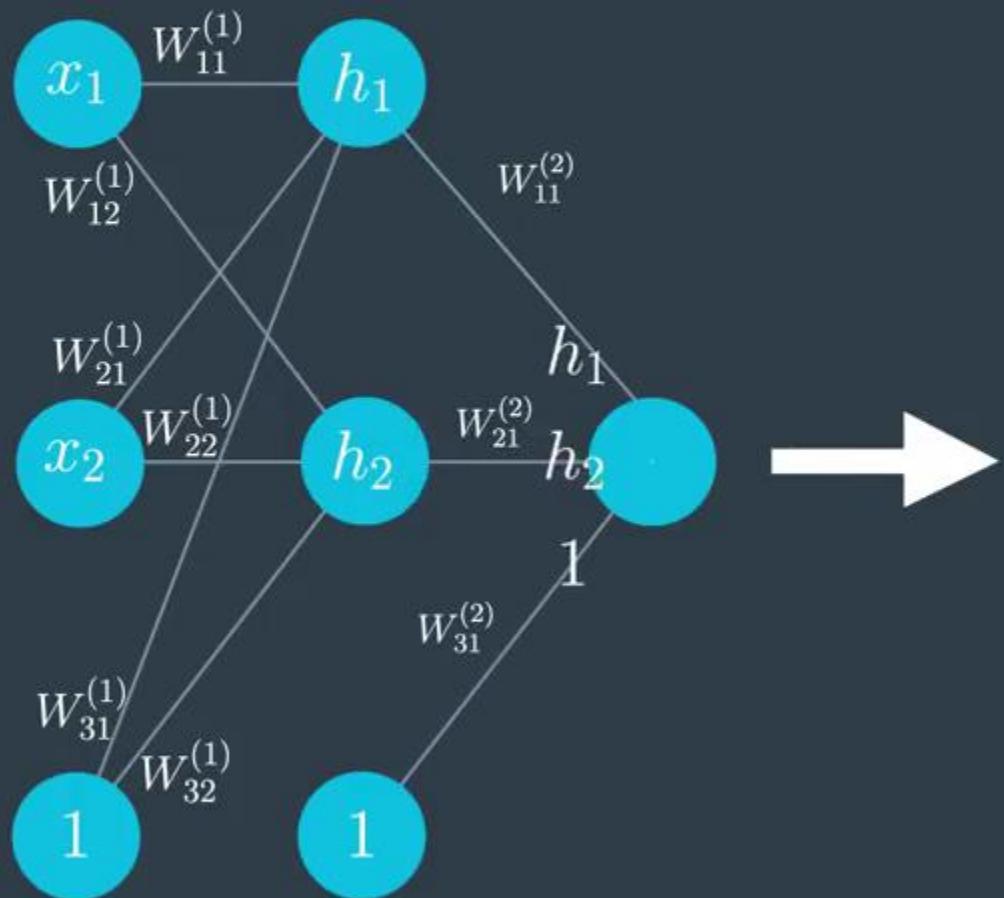


$$h_1 = W_{11}^{(1)}x_1 + W_{21}^{(1)}x_2 + W_{31}^{(1)}$$

$$h_2 = W_{12}^{(1)}x_1 + W_{22}^{(1)}x_2 + W_{32}^{(1)}$$

$$h = W_{11}^{(2)}\sigma(h_1) + W_{21}^{(2)}\sigma(h_2) + W_{31}^{(2)}$$

Feedforward

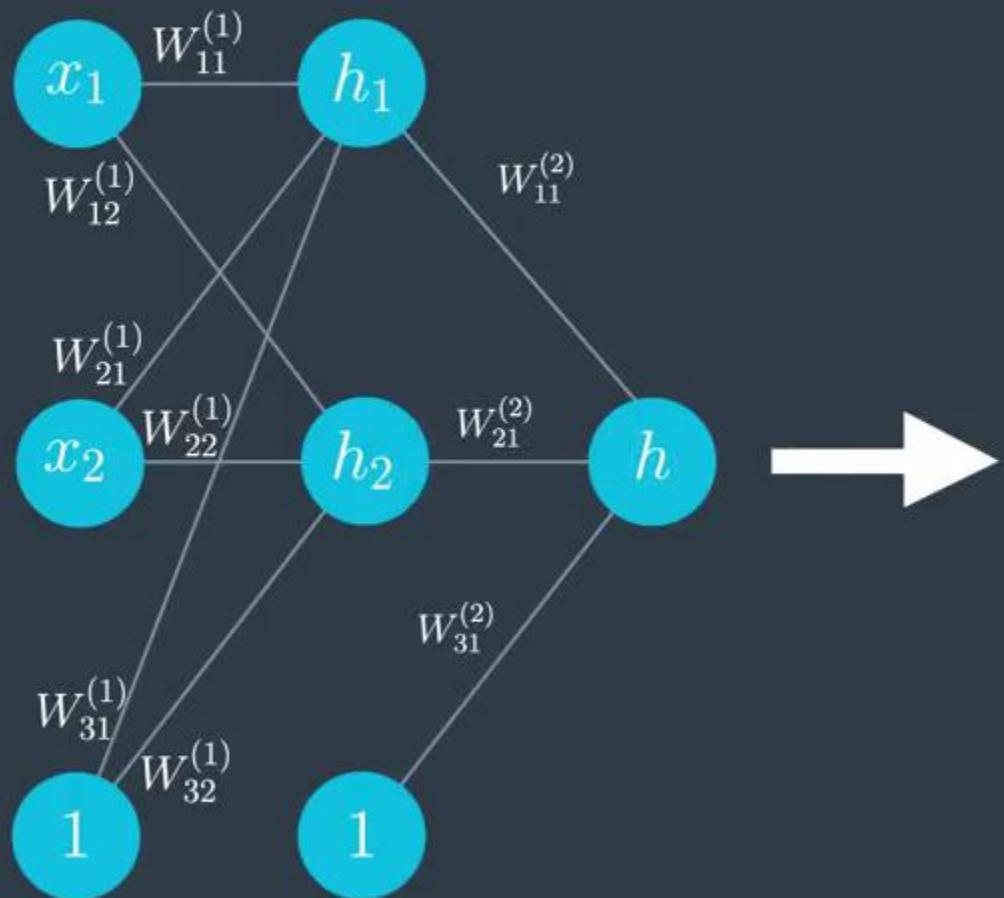


$$h_1 = W_{11}^{(1)}x_1 + W_{21}^{(1)}x_2 + W_{31}^{(1)}$$

$$h_2 = W_{12}^{(1)}x_1 + W_{22}^{(1)}x_2 + W_{32}^{(1)}$$

$$h = W_{11}^{(2)}\sigma(h_1) + W_{21}^{(2)}\sigma(h_2) + W_{31}^{(2)}$$

Feedforward

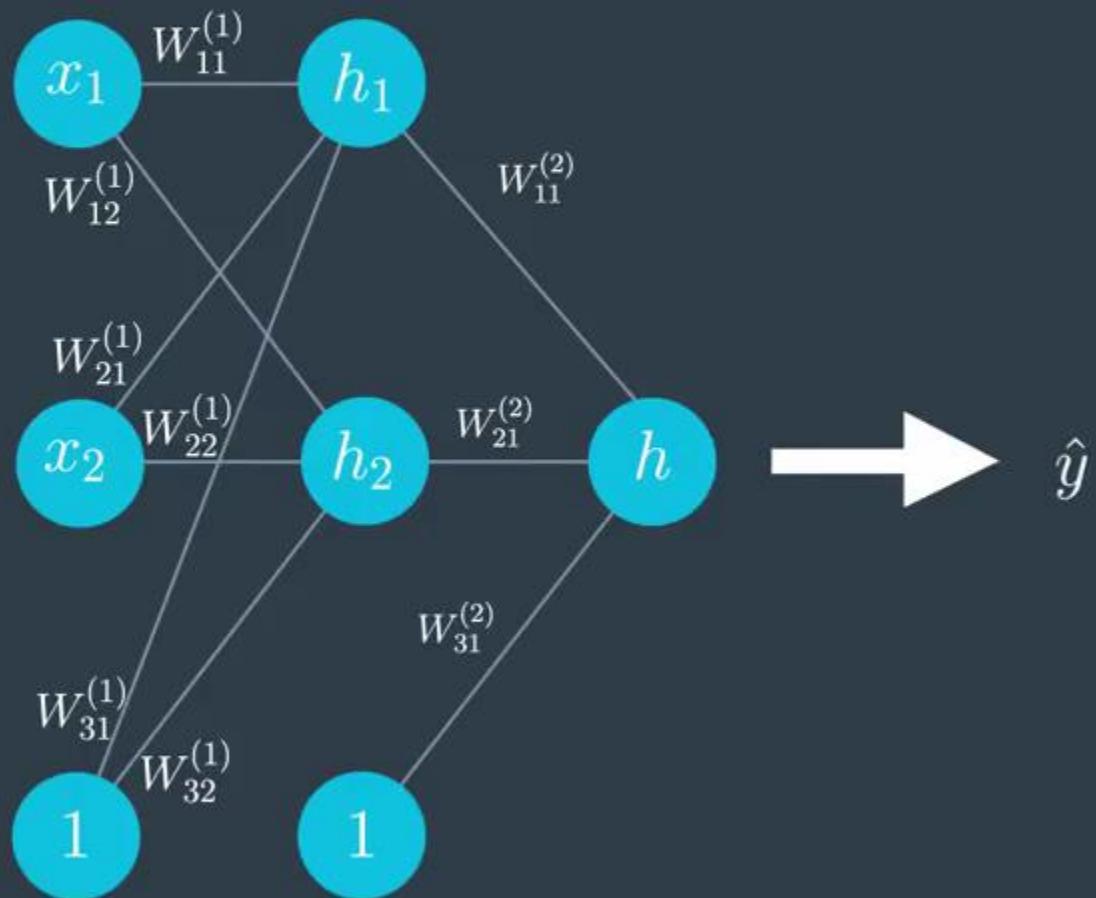


$$h_1 = W_{11}^{(1)}x_1 + W_{21}^{(1)}x_2 + W_{31}^{(1)}$$

$$h_2 = W_{12}^{(1)}x_1 + W_{22}^{(1)}x_2 + W_{32}^{(1)}$$

$$h = W_{11}^{(2)}\sigma(h_1) + W_{21}^{(2)}\sigma(h_2) + W_{31}^{(2)}$$

Feedforward



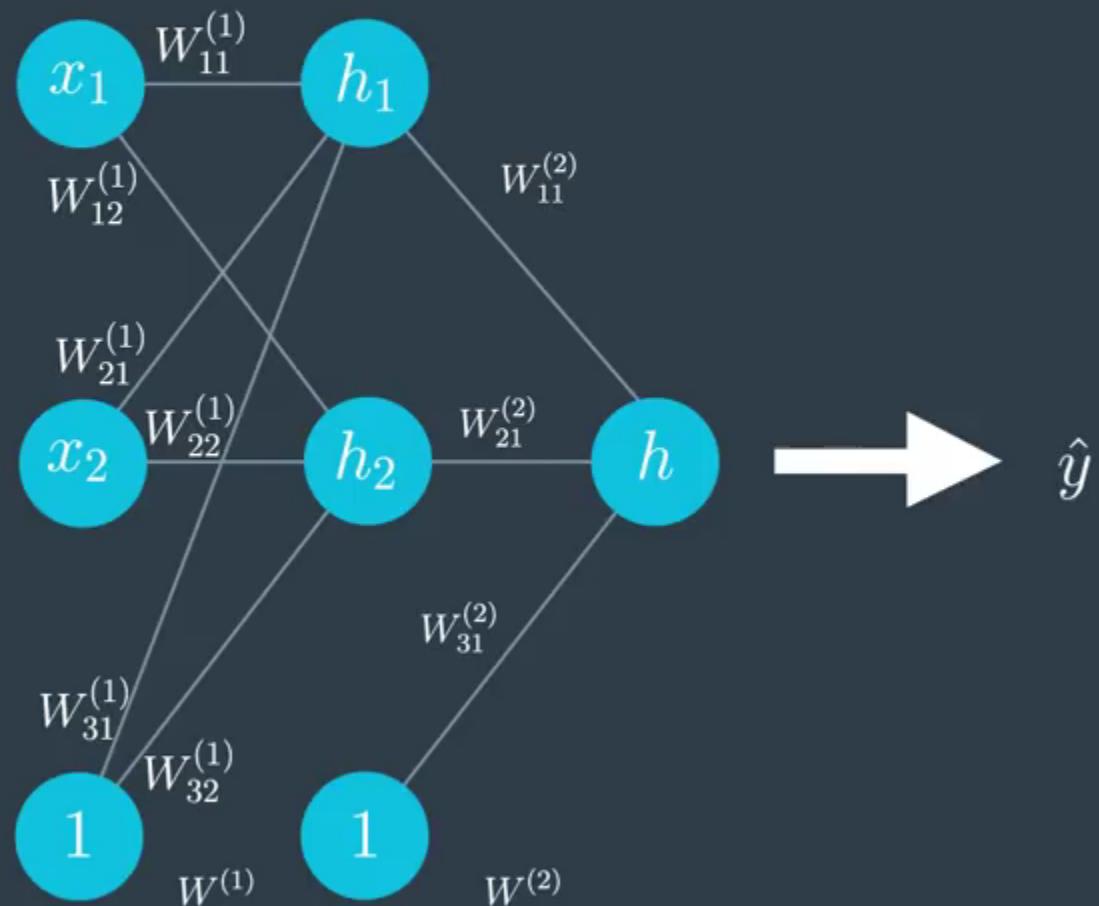
$$h_1 = W_{11}^{(1)}x_1 + W_{21}^{(1)}x_2 + W_{31}^{(1)}$$

$$h_2 = W_{12}^{(1)}x_1 + W_{22}^{(1)}x_2 + W_{32}^{(1)}$$

$$h = W_{11}^{(2)}\sigma(h_1) + W_{21}^{(2)}\sigma(h_2) + W_{31}^{(2)}$$

$$\hat{y} = \sigma(h)$$

Feedforward



$$h_1 = W_{11}^{(1)}x_1 + W_{21}^{(1)}x_2 + W_{31}^{(1)}$$

$$h_2 = W_{12}^{(1)}x_1 + W_{22}^{(1)}x_2 + W_{32}^{(1)}$$

$$h = W_{11}^{(2)}\sigma(h_1) + W_{21}^{(2)}\sigma(h_2) + W_{31}^{(2)}$$

$$\hat{y} = \sigma(h)$$

$$\hat{y} = \sigma \circ W^{(2)} \circ \sigma \circ W^{(1)}(x)$$

Sigmoid Derivative

$$\begin{aligned}\sigma'(x) &= \frac{\partial}{\partial x} \frac{1}{1+e^{-x}} \\&= \frac{e^{-x}}{(1+e^{-x})^2} \\&= \frac{1}{1+e^{-x}} \cdot \frac{e^{-x}}{1+e^{-x}} \\&= \sigma(x)(1 - \sigma(x)) \\&= \sigma(\tilde{x})(1 - \sigma(\tilde{x}))\end{aligned}$$

$$\frac{dh}{dz} = w(1-w)$$

$$L(h, y) = \underbrace{-[y \log(h) + (1-y) \log(1-h)]}_{\text{Log loss derivative}}$$

$$\frac{\partial (C \log(x))}{\partial x} = \frac{C}{x}$$

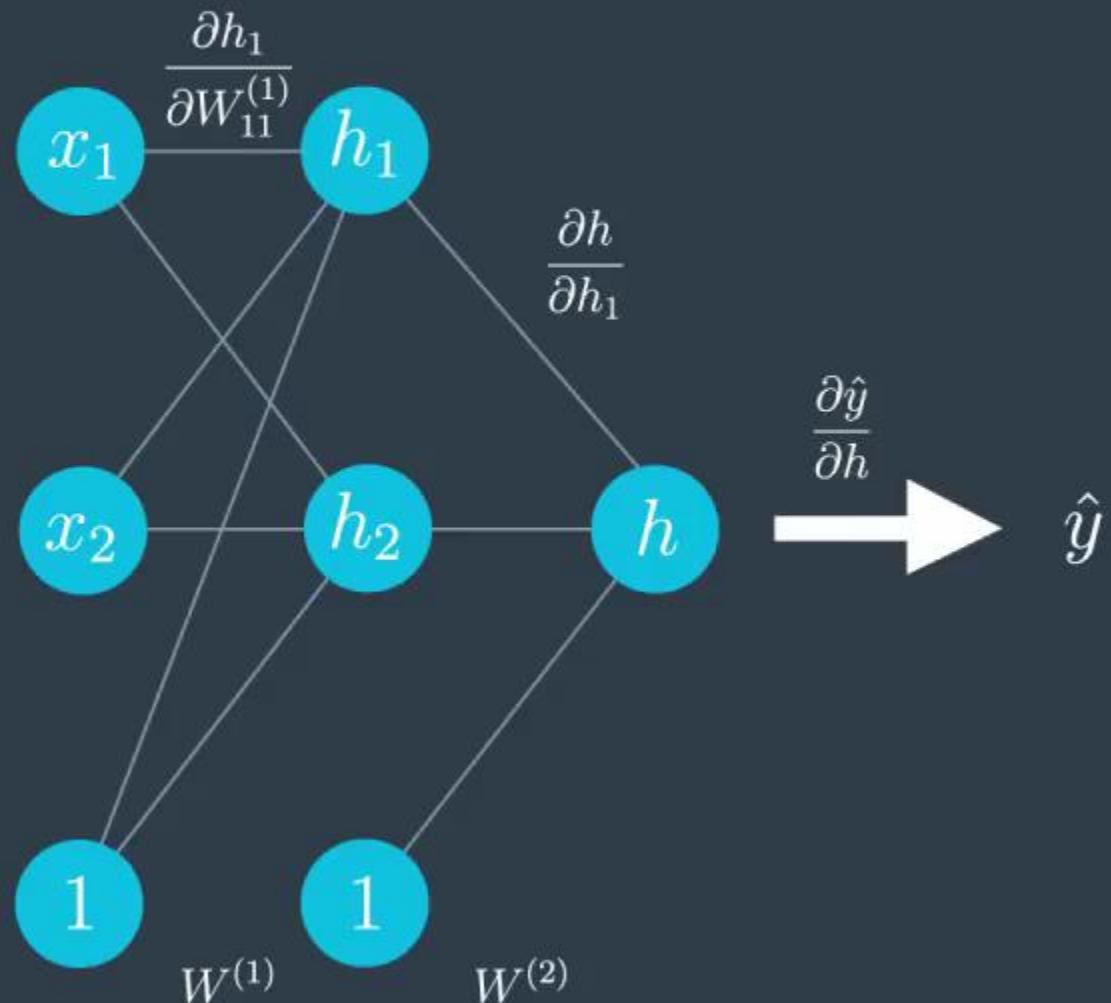
$$\therefore L(h, y) = -\frac{y}{h} - \frac{1-y}{1-h}$$

By using Chain Rule

$$\frac{dL}{dz} = \frac{dL}{dh} - \frac{dh}{dz}$$

$$\begin{aligned}\frac{dL}{dz} &= \left(-\frac{y}{h} - \frac{1-y}{1-h} \right) (h(1-h)) \\ &= -y(1-h) - (1-y)h\end{aligned}$$

BACKPROPAGATION

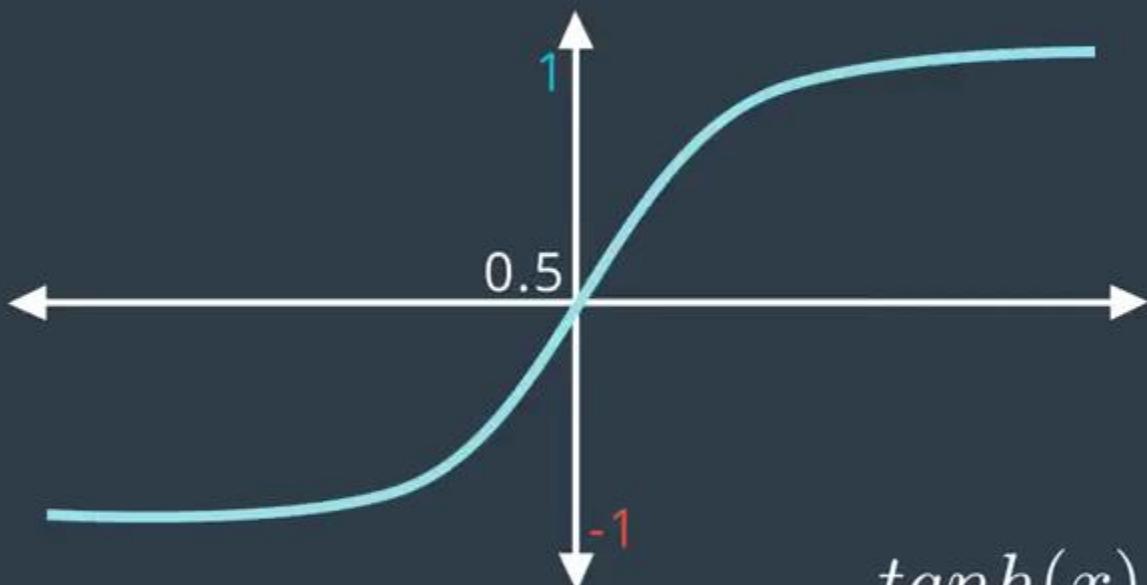


$$\boxed{\frac{\partial E}{\partial W_{11}^{(1)}}} = \boxed{\frac{\partial E}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial h} \frac{\partial h}{\partial h_1} \frac{\partial h_1}{\partial W_{11}^{(1)}}}$$

TINY

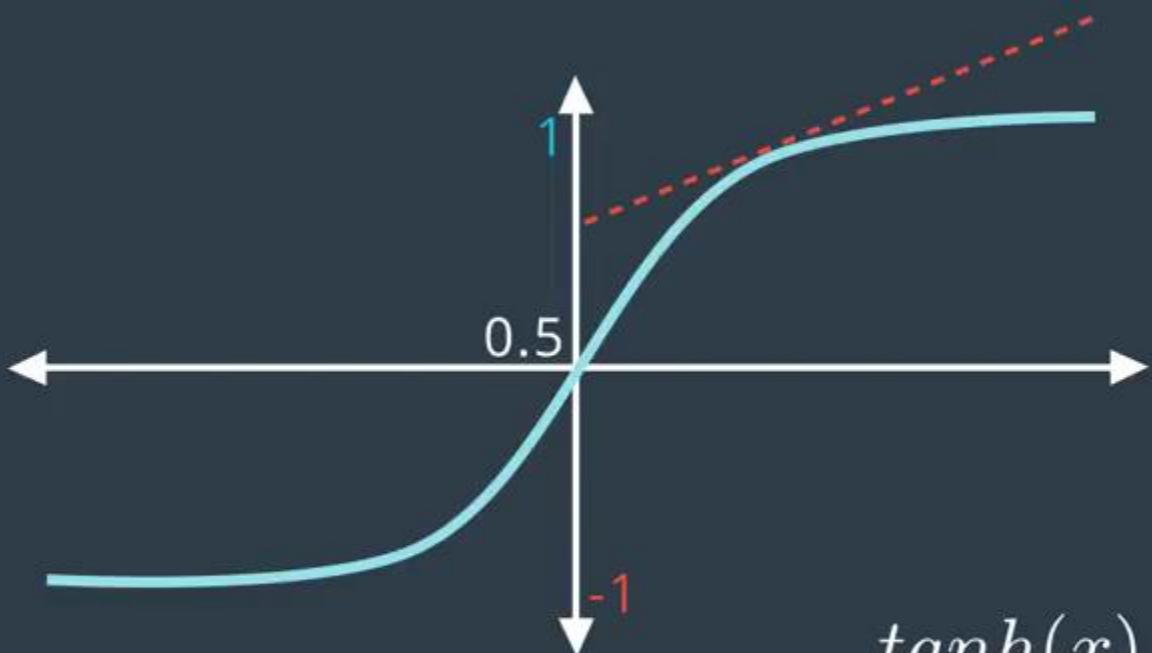
SMALL

HYPERBOLIC TANGENT FUNCTION



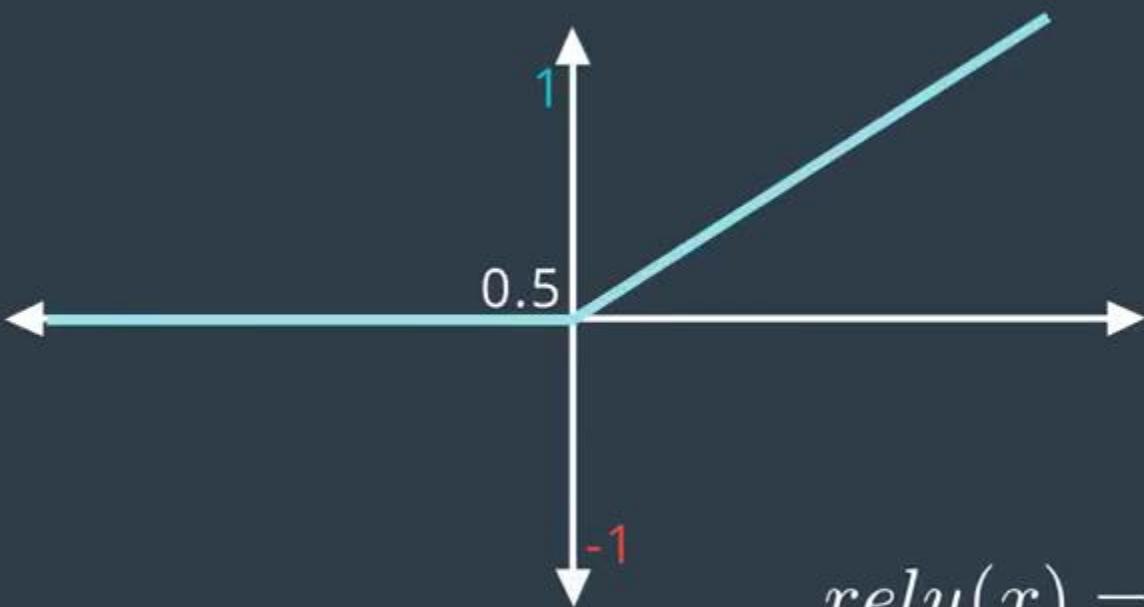
$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

HYPERBOLIC TANGENT FUNCTION



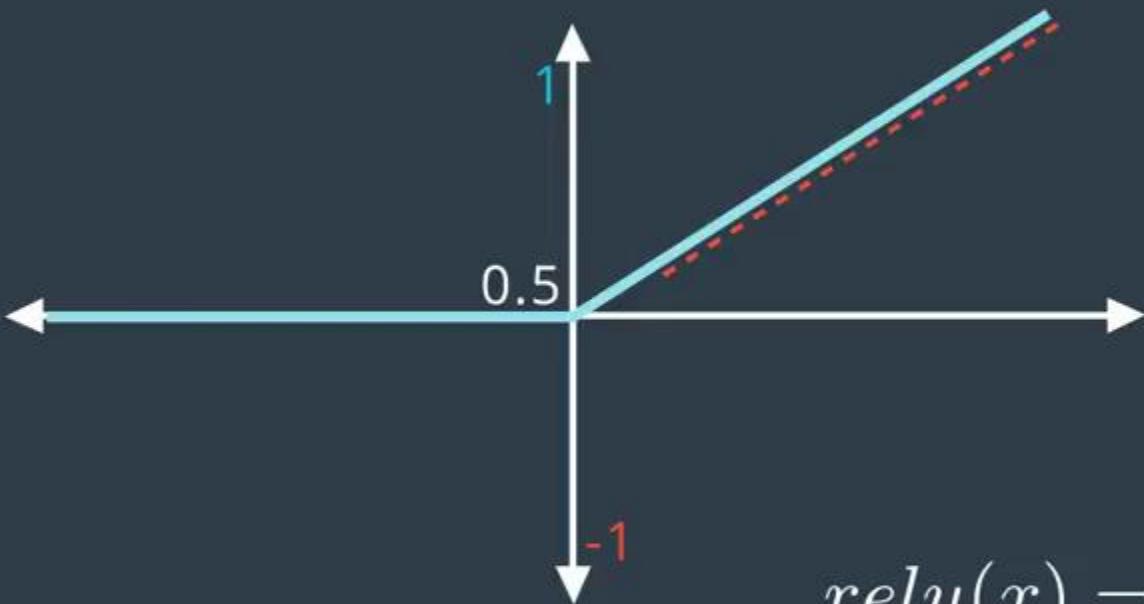
$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

RECTIFIED LINEAR UNIT (ReLU)



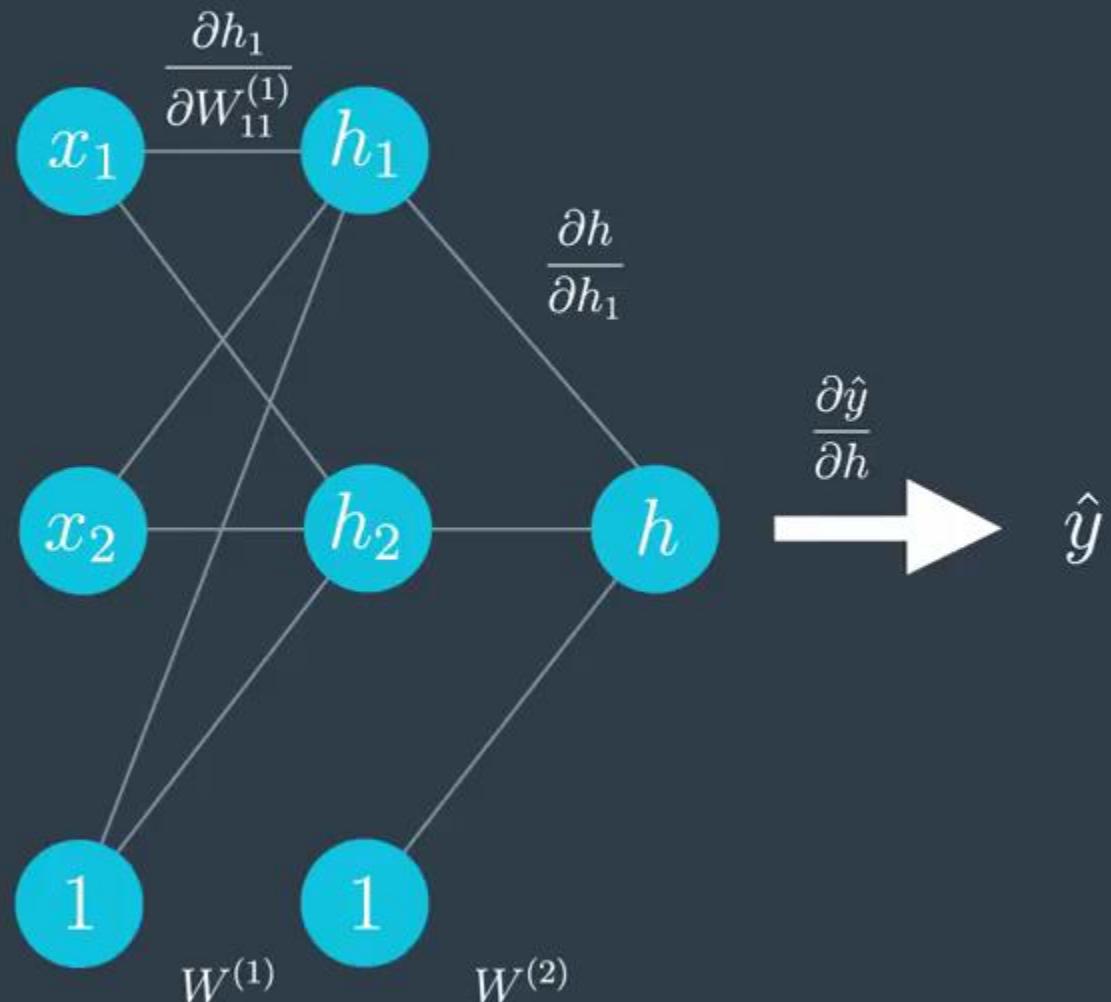
$$\text{relu}(x) = \begin{cases} x & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases}$$

RECTIFIED LINEAR UNIT (ReLU)



$$relu(x) = \begin{cases} x & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases}$$

BACKPROPAGATION



$$\frac{\partial E}{\partial W_{11}^{(1)}} = \frac{\partial E}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial h} \frac{\partial h}{\partial h_1} \frac{\partial h_1}{\partial W_{11}^{(1)}}$$

OK!

OK