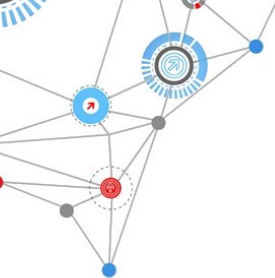


Claudio Mirabello

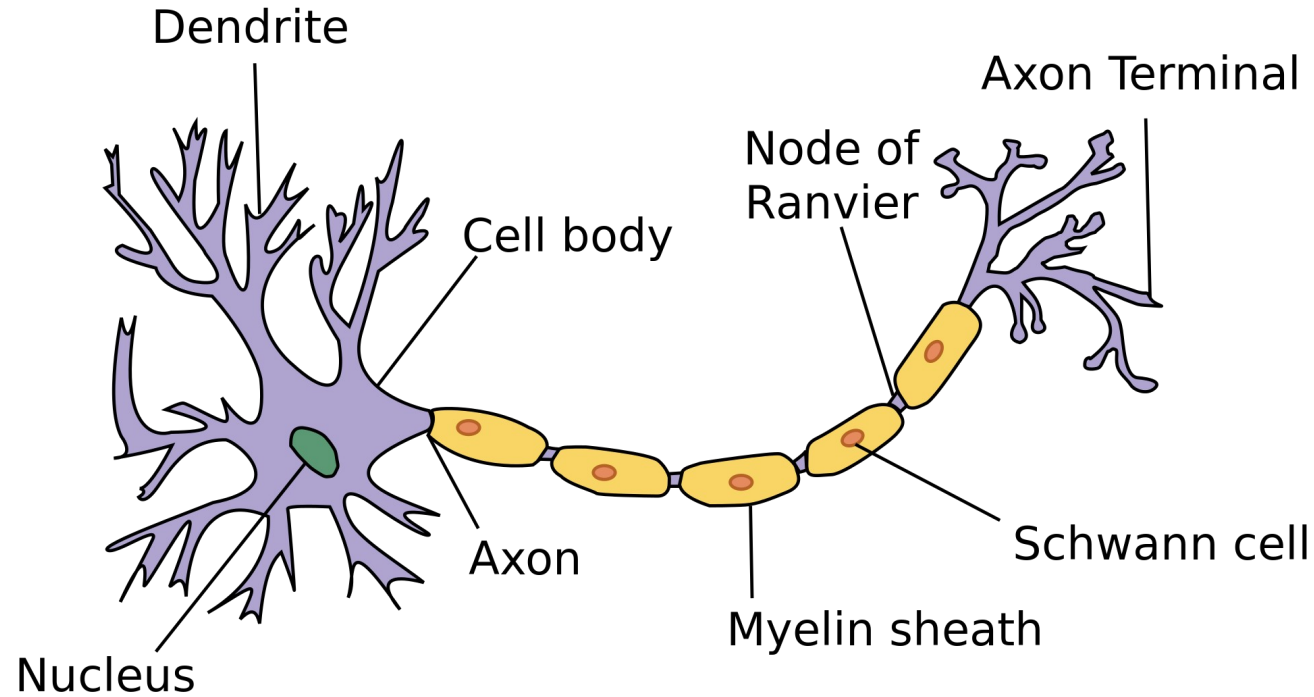
# Basic theory of Feed-forward Neural Networks



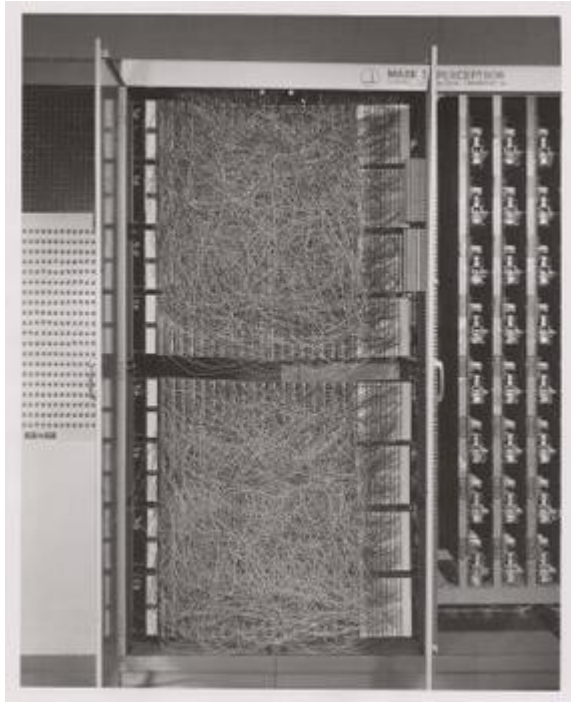
# Resources

- ▶ MIT lectures on Deep Learning  
(<http://introtodeeplearning.com/>)
- ▶ TensorFlow Playground  
(<https://playground.tensorflow.org>)
- ▶ Keras Docs  
(<https://keras.io>)

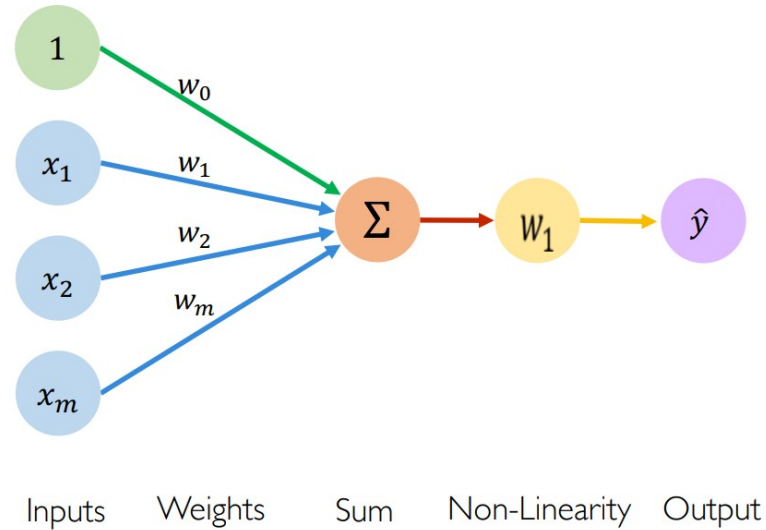
# This is a neuron



# This is a perceptron (1958)



Mark I Perceptron machine  
wikipedia.org



Output

Linear combination of inputs

$$\hat{y} = g \left( w_0 + \sum_{i=1}^m x_i w_i \right)$$

Non-linear activation function

Bias





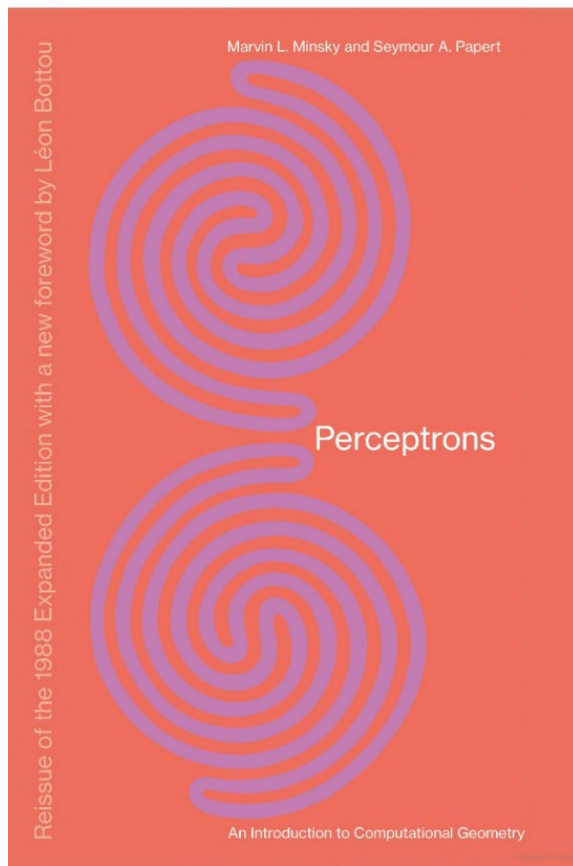
# Perceptrons caused excitement

- ▶ *"the embryo of an electronic computer that [the Navy] expects will be able to walk, talk, see, write, reproduce itself and be conscious of its existence."*

The New York Times



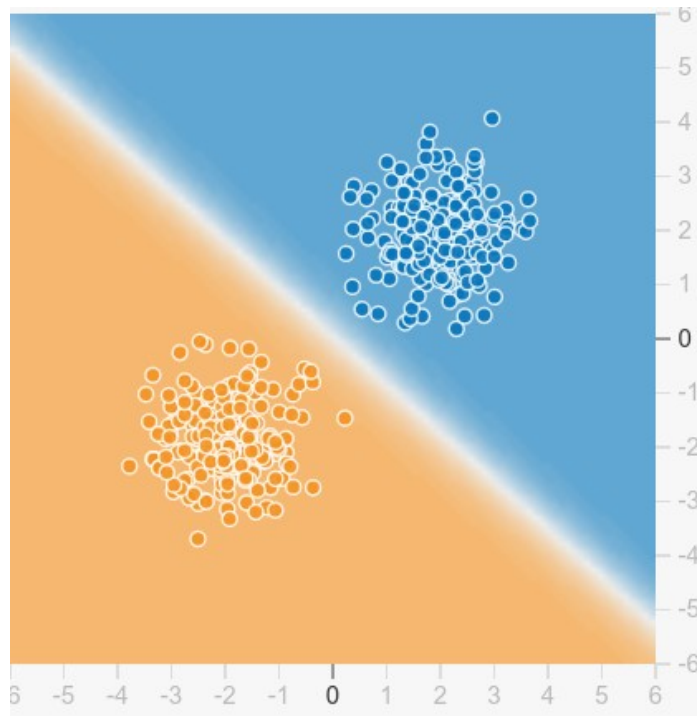
# Perceptrons can only learn linearly separable classes



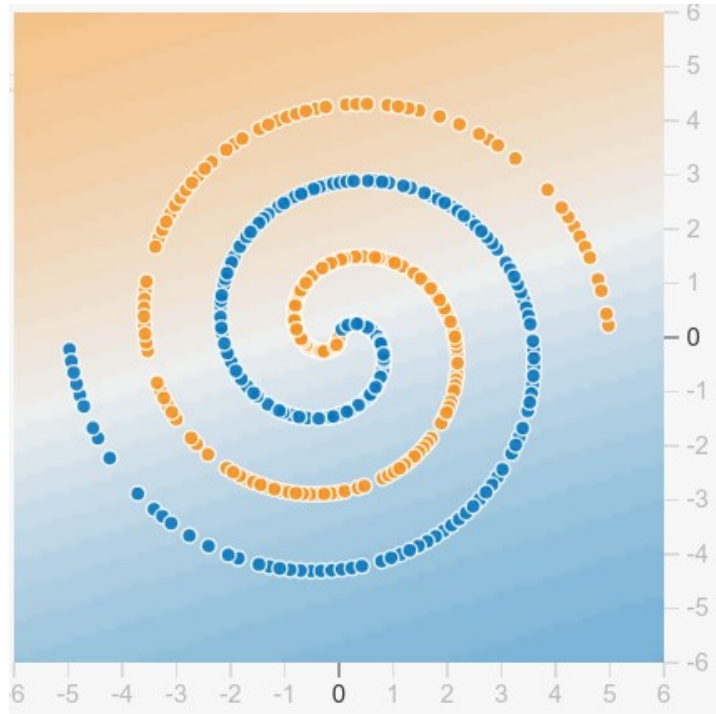
Minsky and Papert, 1969



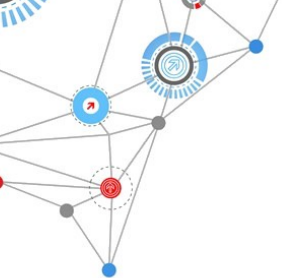
# Perceptrons can only learn linearly separable classes



But sometimes you want  
to model non-linear functions

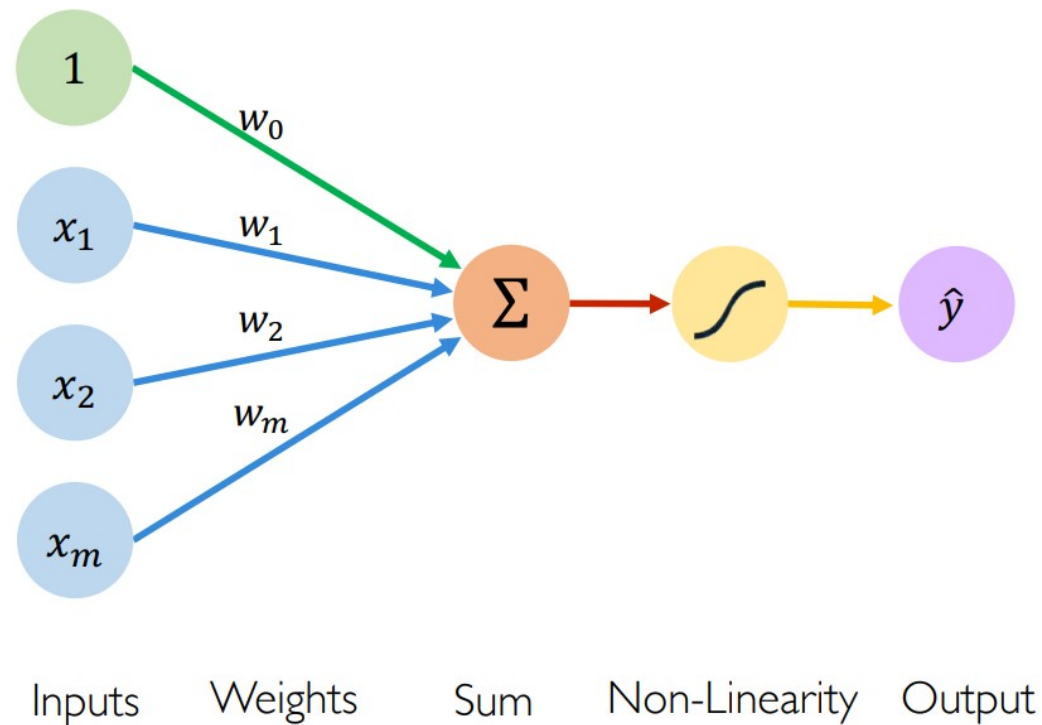






How do we make this non-linear then?  
Two ingredients to add

# 1: Differentiable, non-linear activation functions

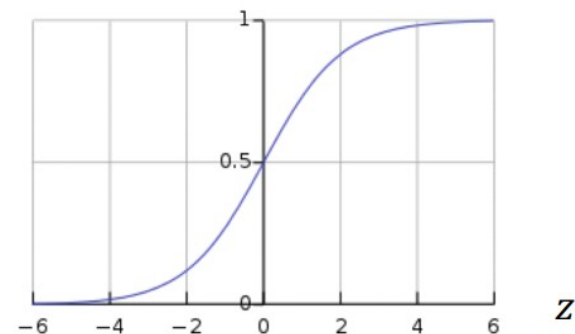


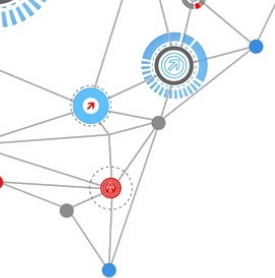
## Activation Functions

$$\hat{y} = g(w_0 + \mathbf{X}^T \mathbf{W})$$

- Example: sigmoid function

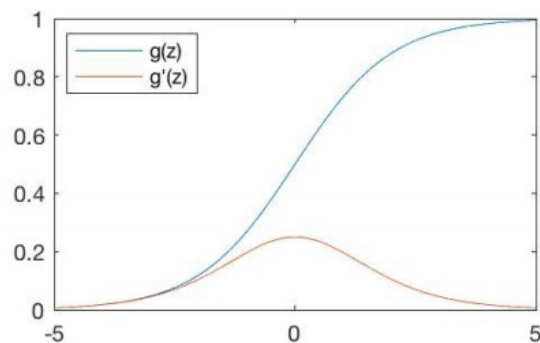
$$g(z) = \sigma(z) = \frac{1}{1 + e^{-z}}$$





# Common activation functions

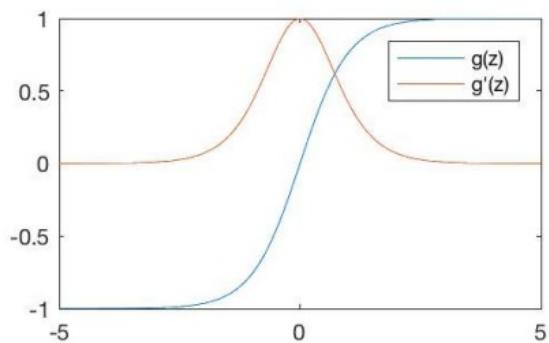
Sigmoid Function



$$g(z) = \frac{1}{1 + e^{-z}}$$

$$g'(z) = g(z)(1 - g(z))$$

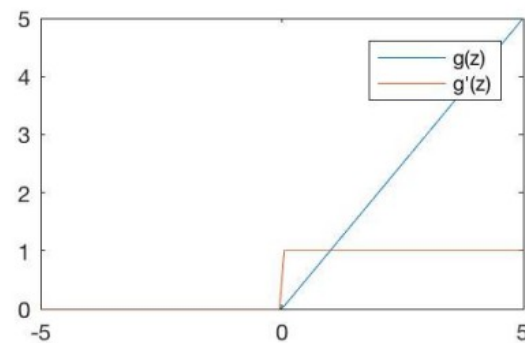
Hyperbolic Tangent



$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$g'(z) = 1 - g(z)^2$$

Rectified Linear Unit (ReLU)



$$g(z) = \max(0, z)$$

$$g'(z) = \begin{cases} 1, & z > 0 \\ 0, & \text{otherwise} \end{cases}$$



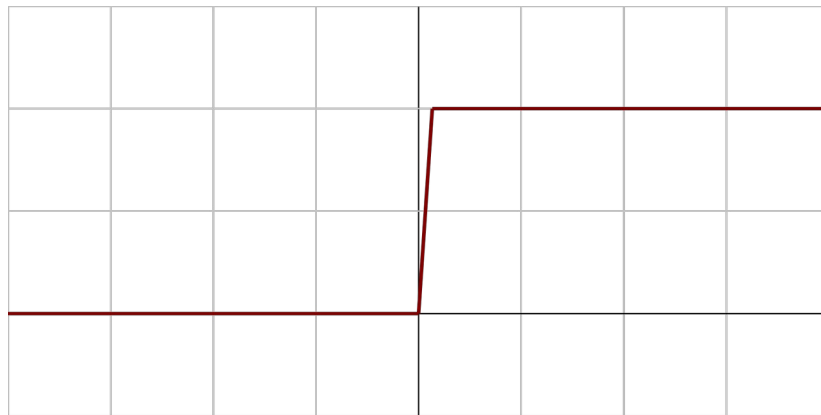
## Special case: softmax

- Used in classification problems
- Given  $k$  classes, it decides which one is more likely
- One output per class, each output is assigned a probability from 0 to 1
- The sum of probabilities for all outputs is 1

$$g(\mathbf{z})_j = \frac{e^{z_j}}{\sum_{k=1}^k e^{z_k}} \text{ for } j = 1, \dots, k$$

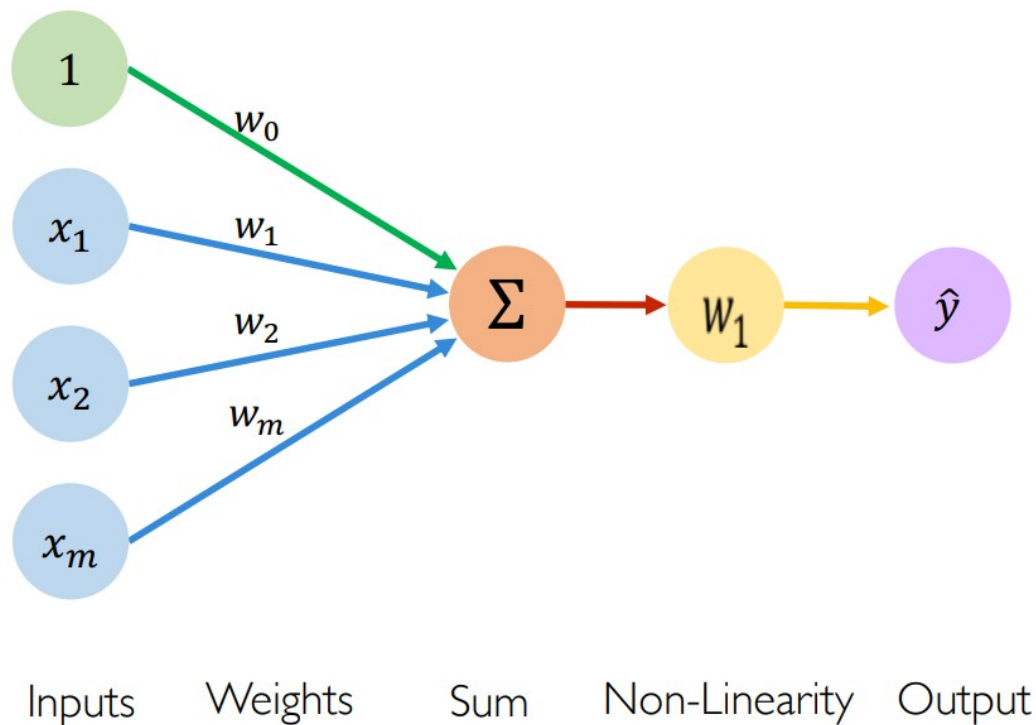


Wait a second, the perceptron already has a non-linear (step) activation function!





# This is a perceptron



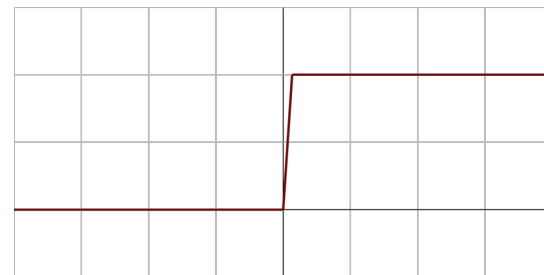
Output

Linear combination of inputs

$$\hat{y} = g \left( w_0 + \sum_{i=1}^m x_i w_i \right)$$

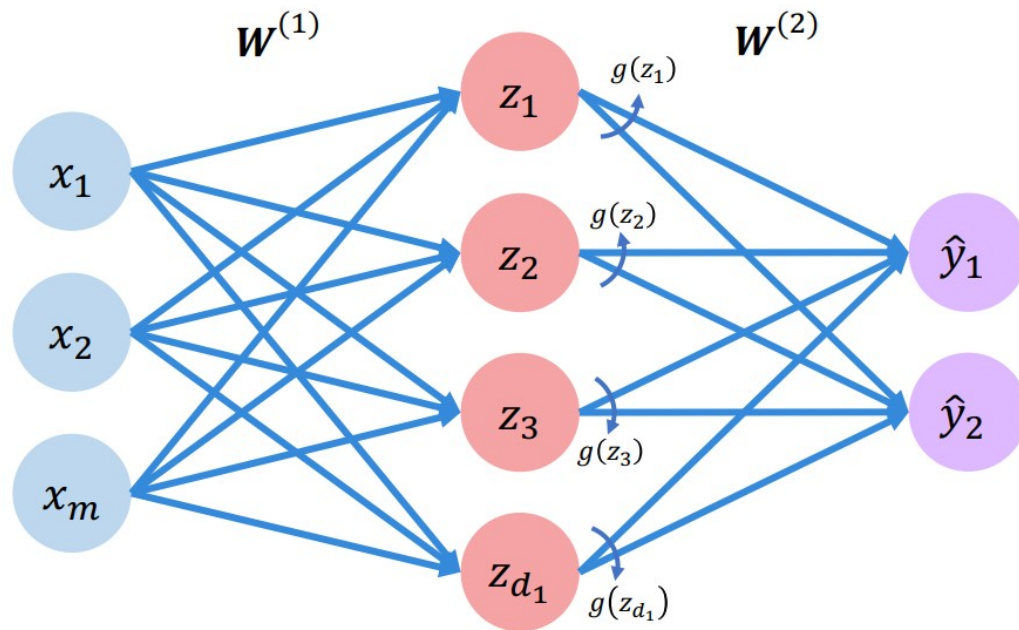
activation function

Bias





## 2: Multi-layer Perceptron (1986)

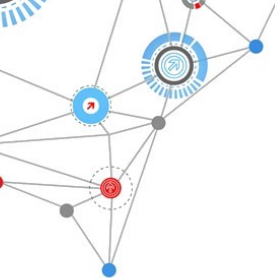


Inputs

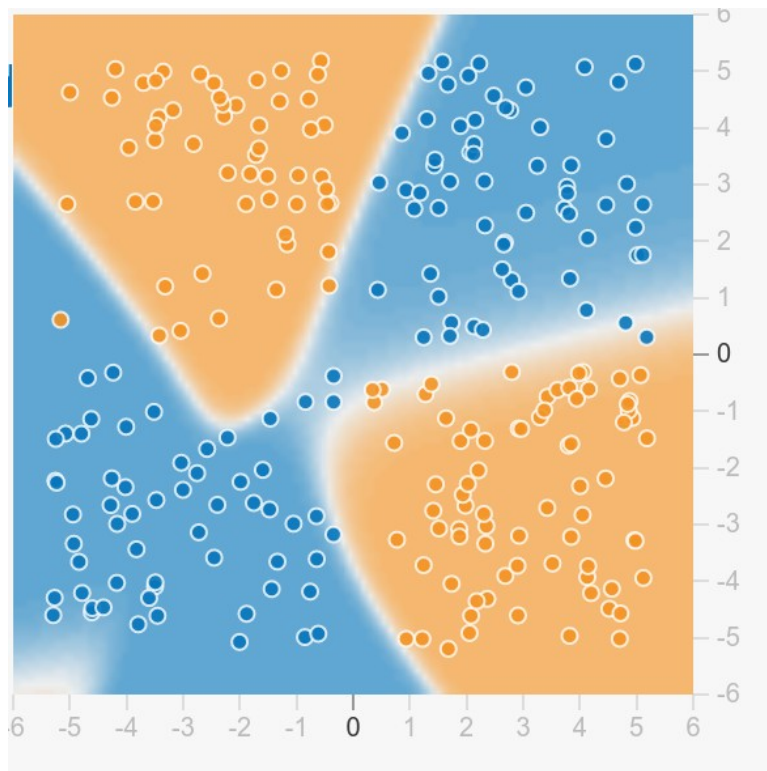
Hidden

Final Output

$$z_i = w_{0,i}^{(1)} + \sum_{j=1}^m x_j w_{j,i}^{(1)} \quad \hat{y}_i = g \left( w_{0,i}^{(2)} + \sum_{j=1}^{d_1} z_j w_{j,i}^{(2)} \right)$$



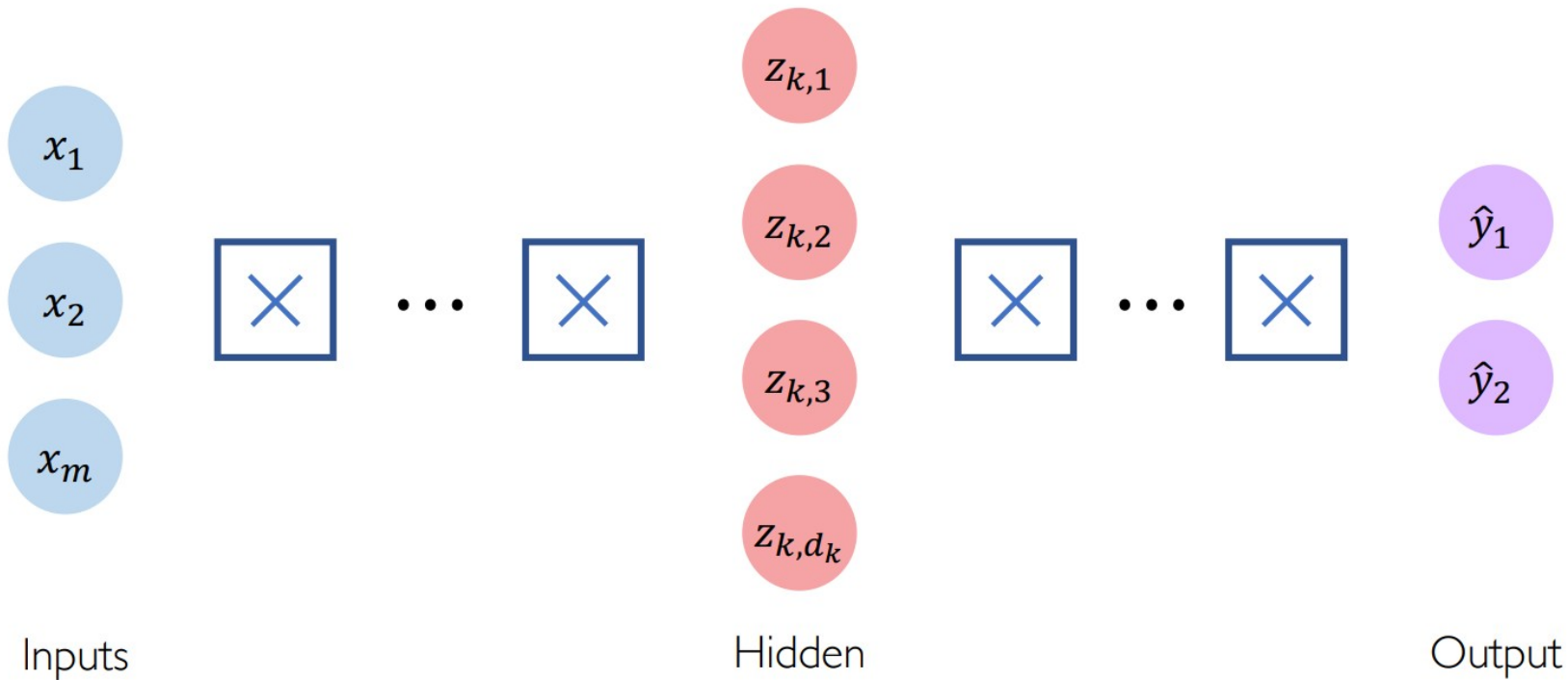
Now we're getting somewhere



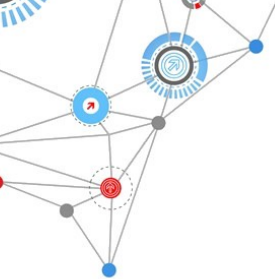




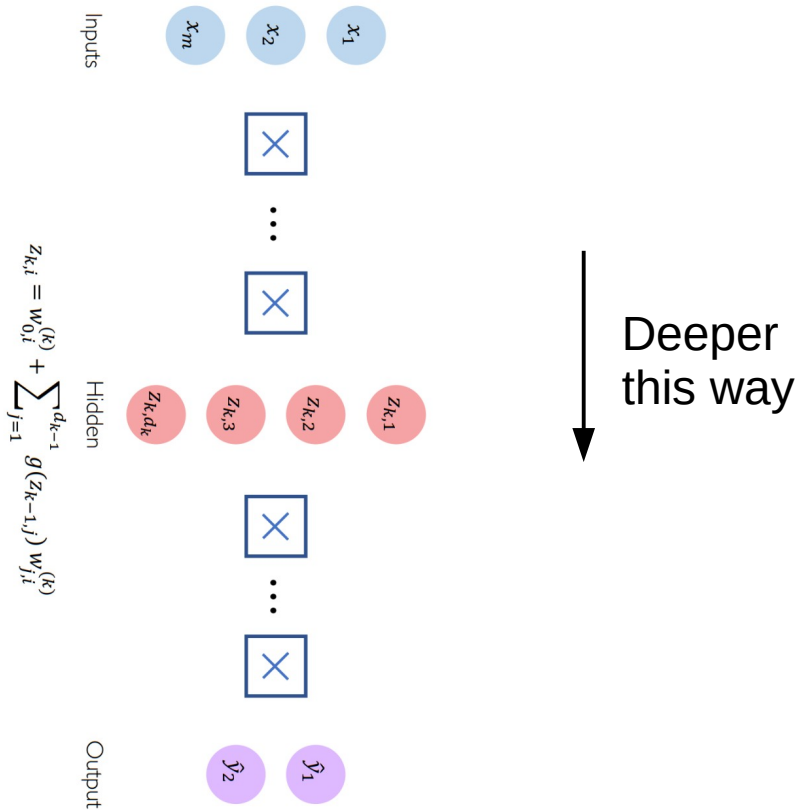
# Why stop at one hidden layer?



$$z_{k,i} = w_{0,i}^{(k)} + \sum_{j=1}^{d_{k-1}} g(z_{k-1,j}) w_{j,i}^{(k)}$$

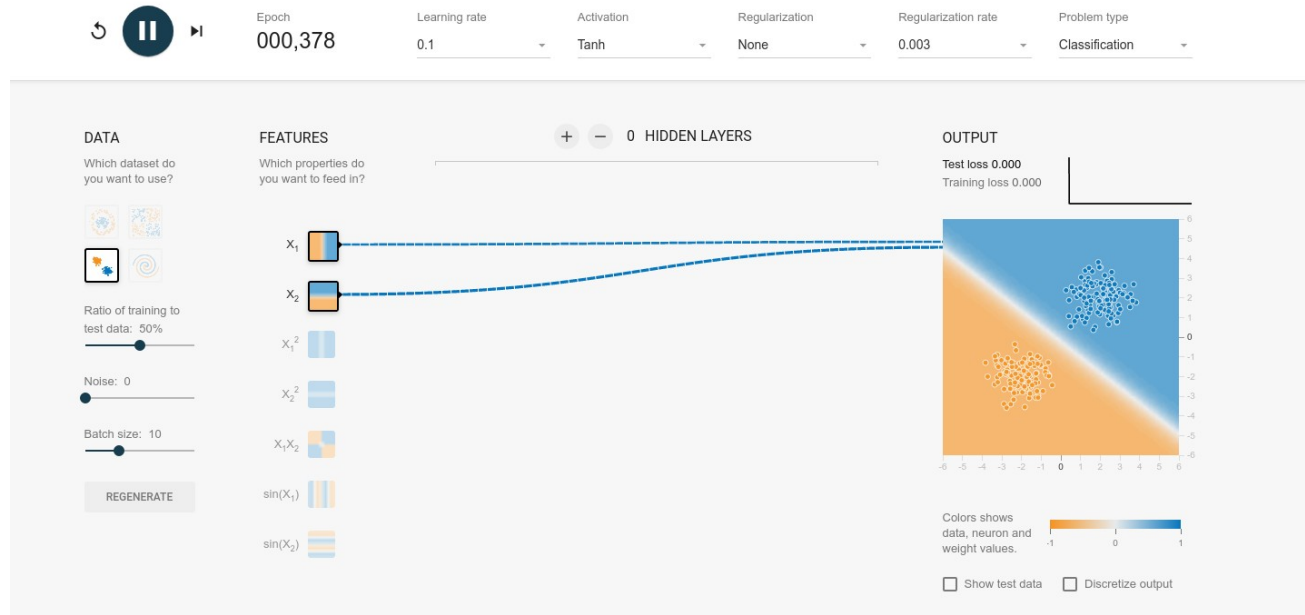


# Deep Networks are simply NNs with multiple hidden layers



# <https://playground.tensorflow.org>

Tinker With a **Neural Network** Right Here in Your Browser.  
Don't Worry, You Can't Break It. We Promise.



Let's review:

- Perceptron
- XOR problem
- Activations
- Multi-layer perceptron



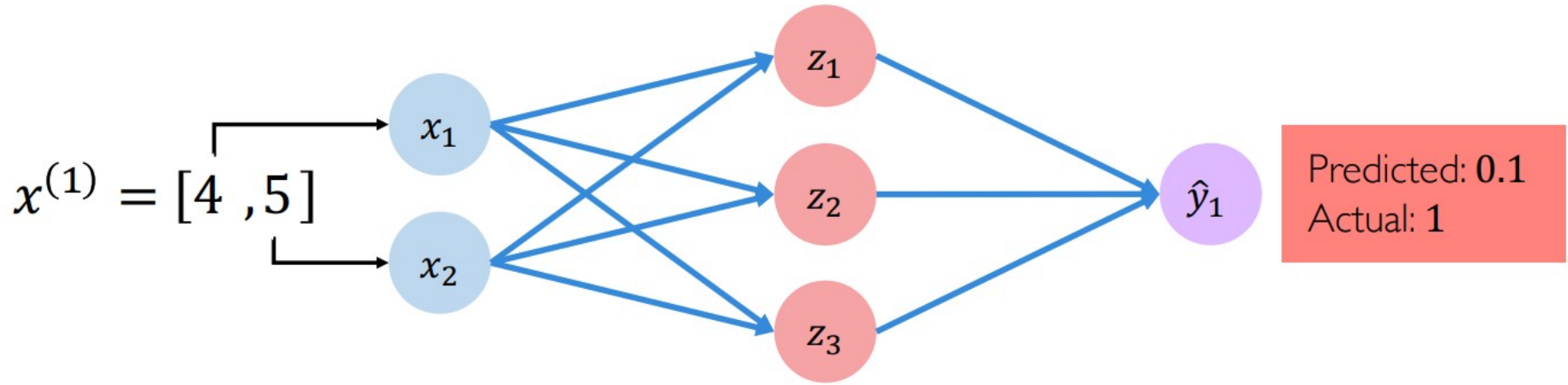
# How do we decide which weights are optimal?

- A linear regressor's weights (coefficients) are calculated in closed form
- This can't be done if you have hidden layers and non-linear activations



# How do we decide which weights are optimal?

The **loss** of our network measures the cost incurred from incorrect predictions

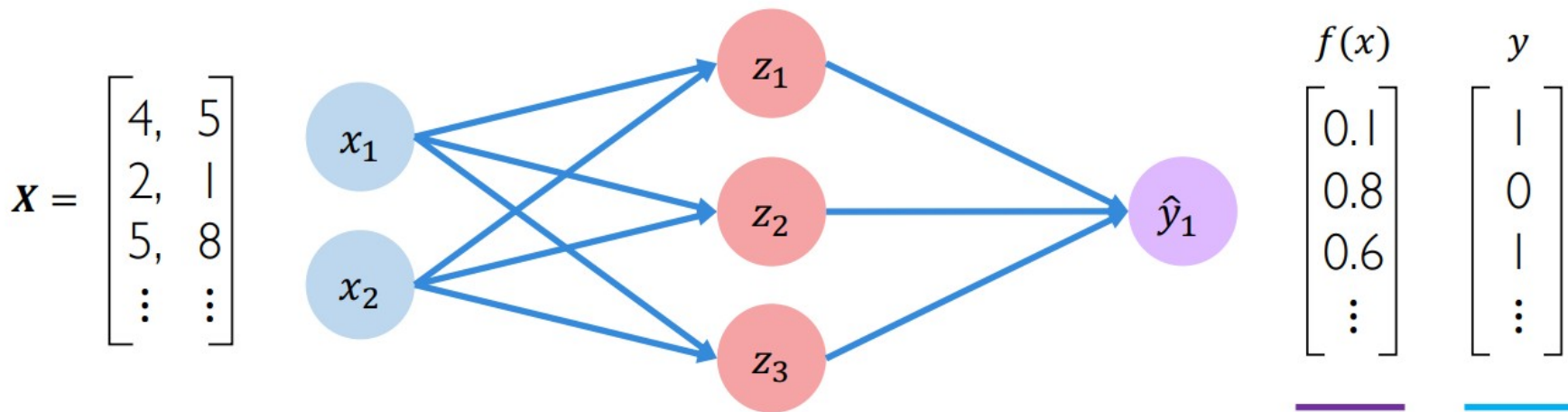


$$\mathcal{L}(\underbrace{f(x^{(i)}; \mathbf{W})}_{\text{Predicted}}, \underbrace{y^{(i)}}_{\text{Actual}})$$



# How do we decide which weights are optimal?

The **empirical loss** measures the total loss over our entire dataset



Also known as:

- Objective function
- Cost function
- Empirical Risk

$$J(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(\underbrace{f(x^{(i)}; \mathbf{W})}_{\text{Predicted}}, \underbrace{y^{(i)}}_{\text{Actual}})$$



# Lower loss => better predictions

We want to find the network weights that **achieve the lowest loss**

$$\mathbf{W}^* = \operatorname{argmin}_{\mathbf{W}} \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x^{(i)}; \mathbf{W}), y^{(i)})$$

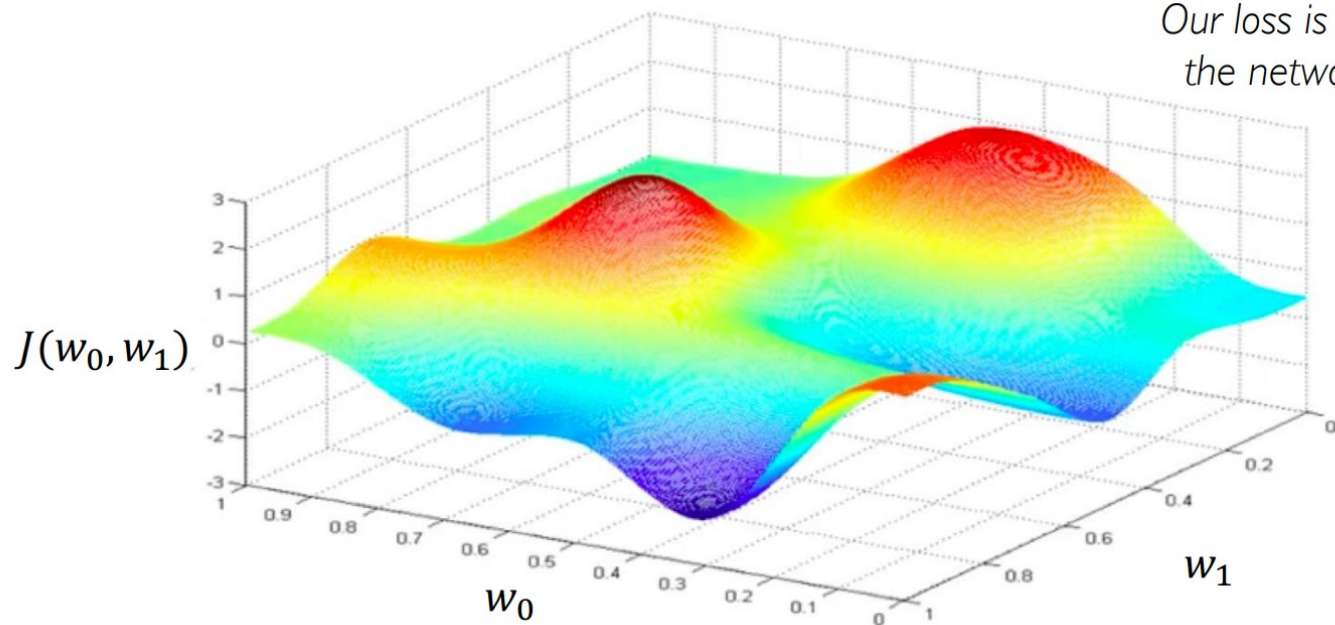
$$\mathbf{W}^* = \operatorname{argmin}_{\mathbf{W}} J(\mathbf{W})$$



# Minimizing loss

$$\mathbf{W}^* = \underset{\mathbf{W}}{\operatorname{argmin}} J(\mathbf{W})$$

Remember:  
*Our loss is a function of  
the network weights!*

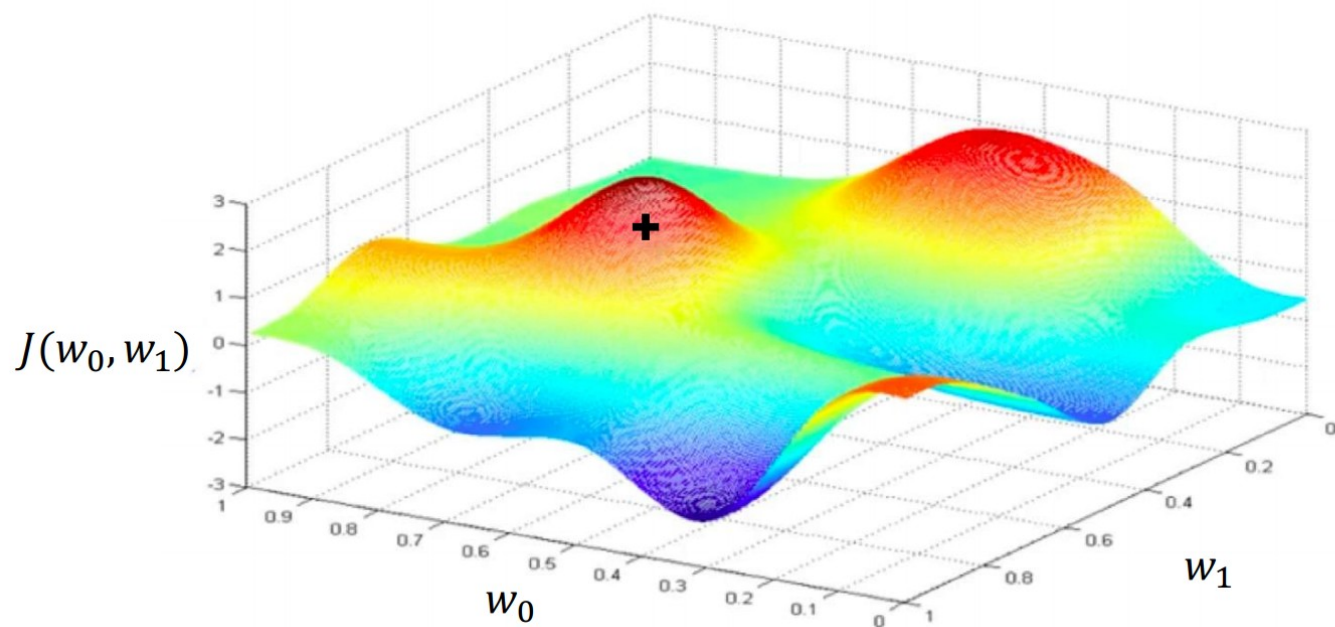


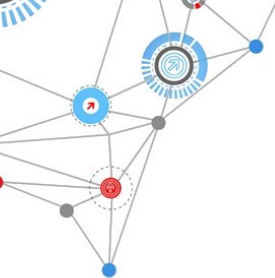




# Minimizing loss

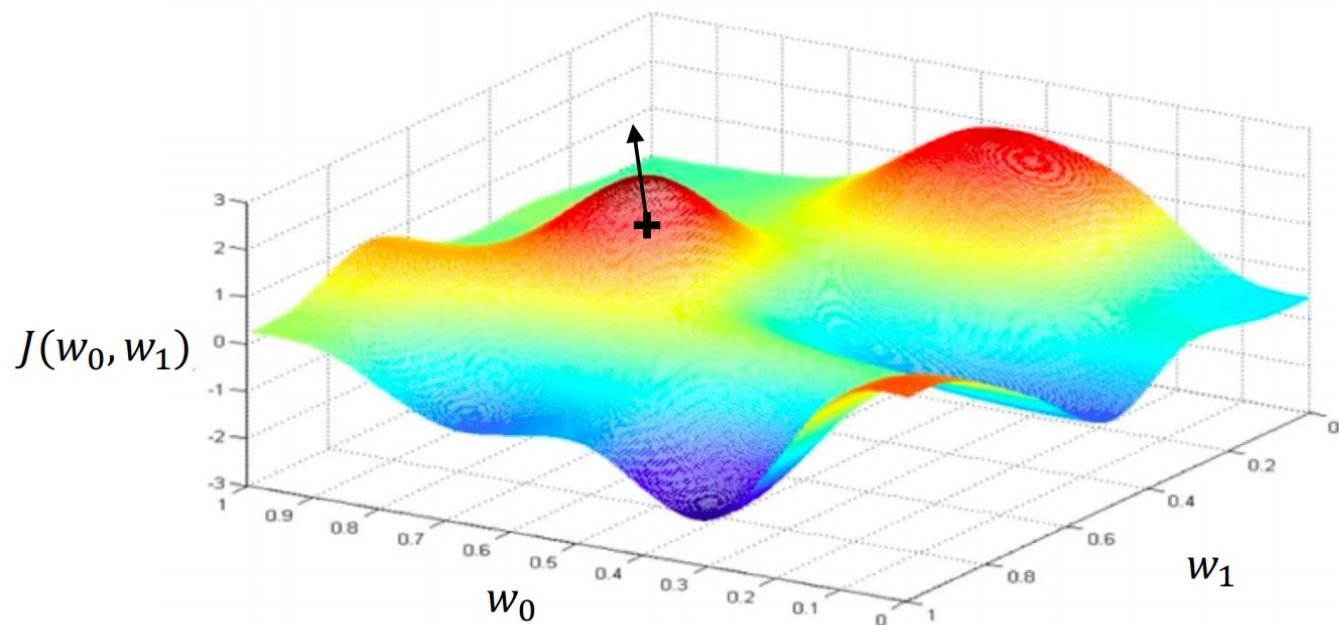
Randomly pick an initial  $(w_0, w_1)$





# Minimizing loss

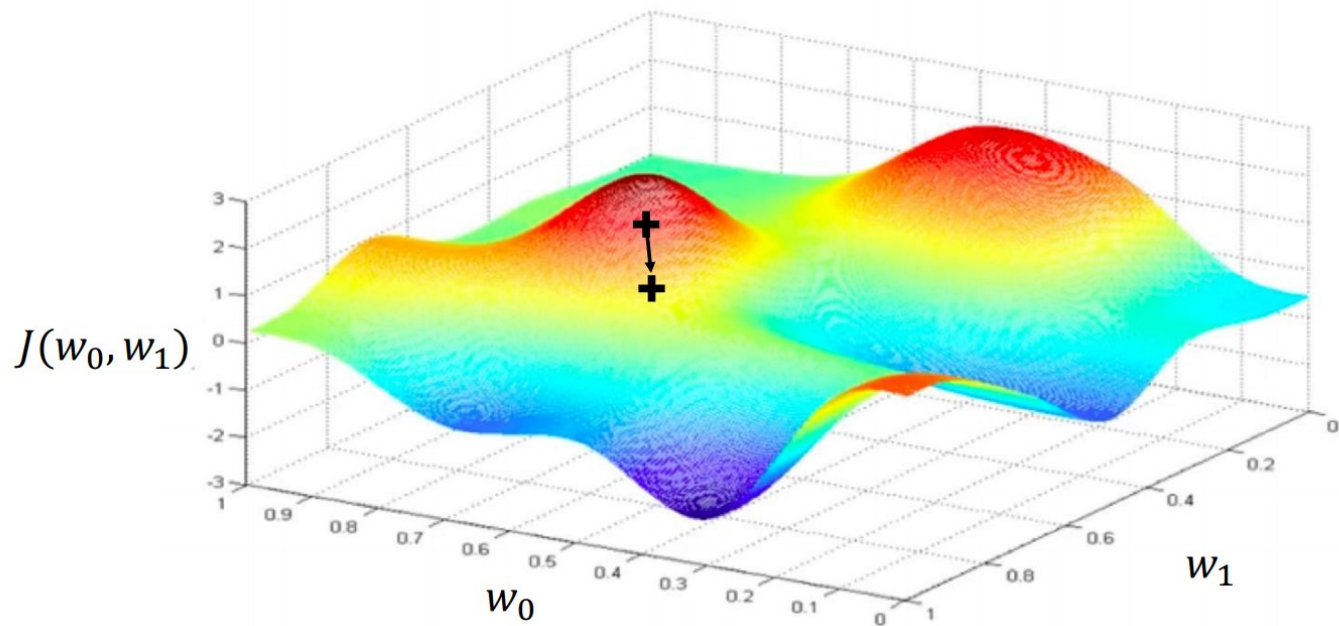
Compute gradient,  $\frac{\partial J(W)}{\partial W}$





# Minimizing loss

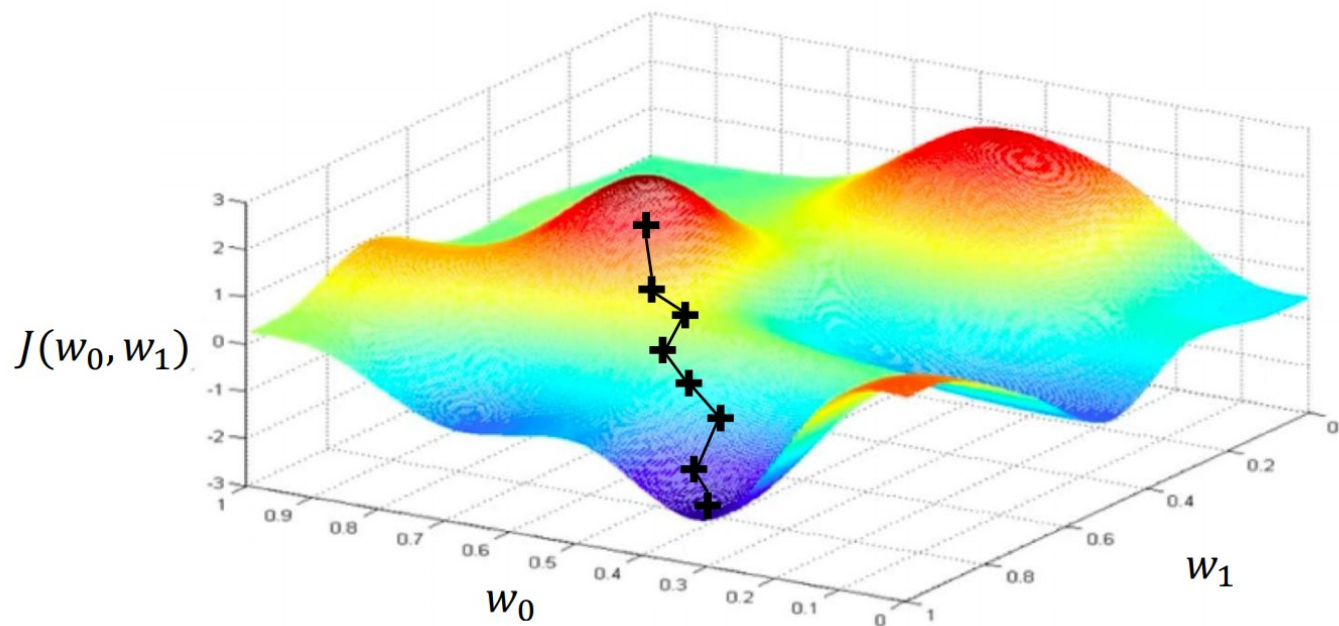
Take small step in opposite direction of gradient





# Minimizing loss

Repeat until convergence





# Gradient descent

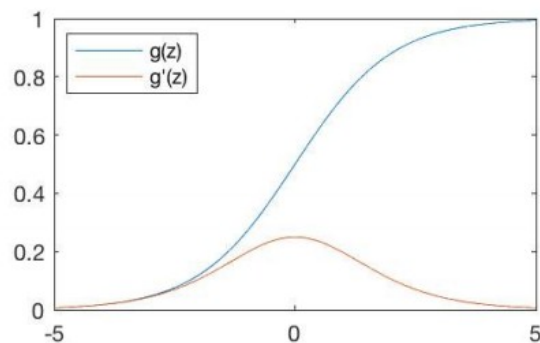
## Algorithm

1. Initialize weights randomly  $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3.     Compute gradient,  $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
4.     Update weights,  $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
5. Return weights



# Activation functions have to be differentiable!

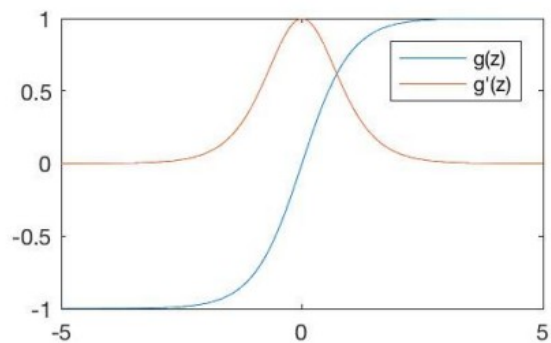
Sigmoid Function



$$g(z) = \frac{1}{1 + e^{-z}}$$

$$g'(z) = g(z)(1 - g(z))$$

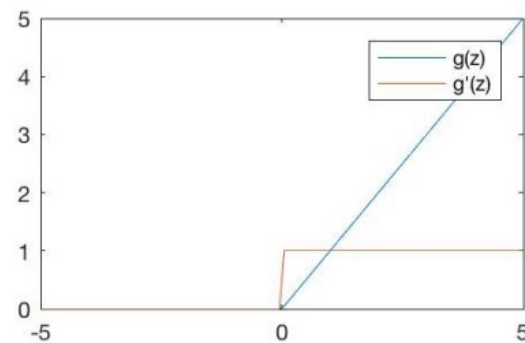
Hyperbolic Tangent



$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$g'(z) = 1 - g(z)^2$$

Rectified Linear Unit (ReLU)



$$g(z) = \max(0, z)$$

$$g'(z) = \begin{cases} 1, & z > 0 \\ 0, & \text{otherwise} \end{cases}$$

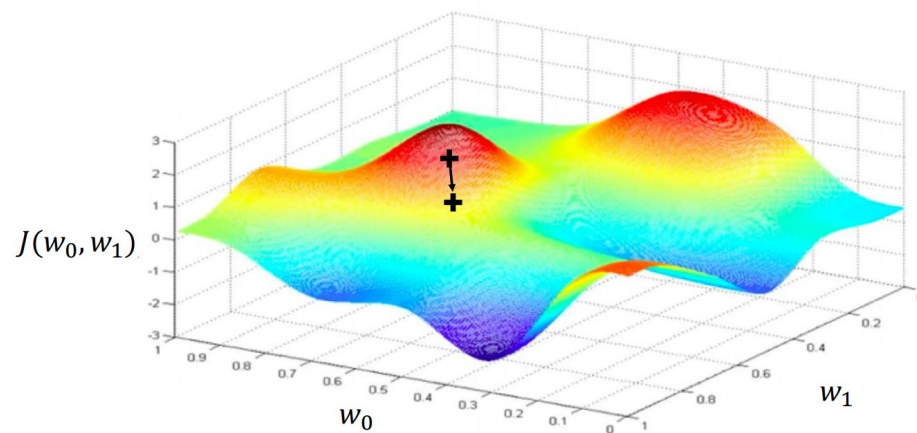


# The learning rate $\eta$

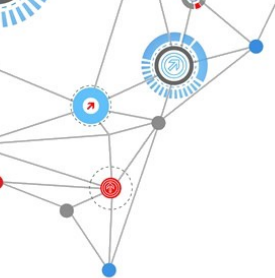
## Algorithm

1. Initialize weights randomly  $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3.     Compute gradient,  $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
4.     Update weights,  $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
5. Return weights

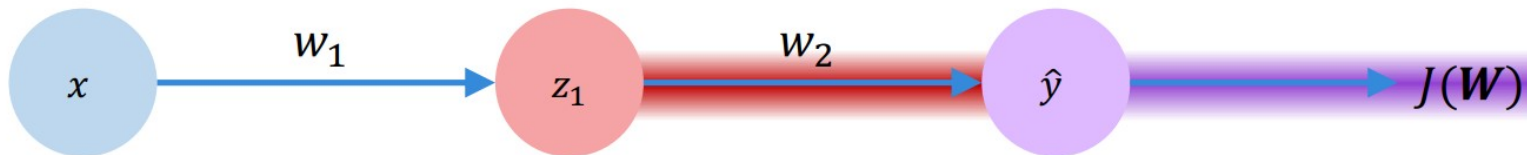
Take small step in opposite direction of gradient







# Backpropagation

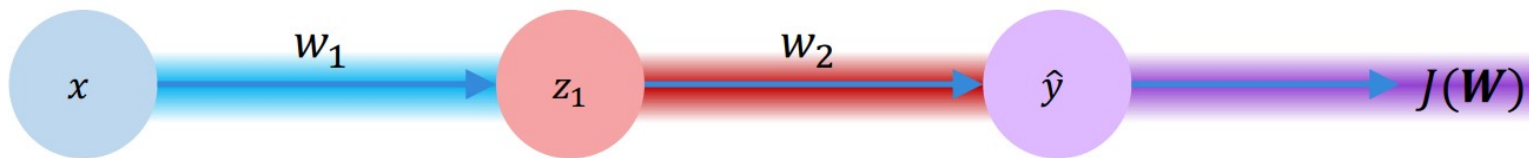


$$\frac{\partial J(W)}{\partial w_2} = \underbrace{\frac{\partial J(W)}{\partial \hat{y}}}_{\text{purple}} * \underbrace{\frac{\partial \hat{y}}{\partial w_2}}_{\text{red}}$$





# Backpropagation



$$\frac{\partial J(\mathbf{W})}{\partial w_1} = \underbrace{\frac{\partial J(\mathbf{W})}{\partial \hat{y}}}_{\text{purple}} * \underbrace{\frac{\partial \hat{y}}{\partial z_1}}_{\text{red}} * \underbrace{\frac{\partial z_1}{\partial w_1}}_{\text{blue}}$$

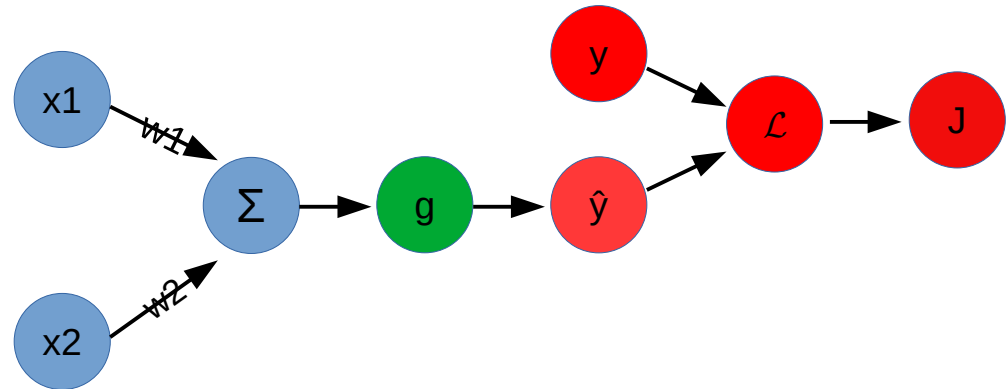
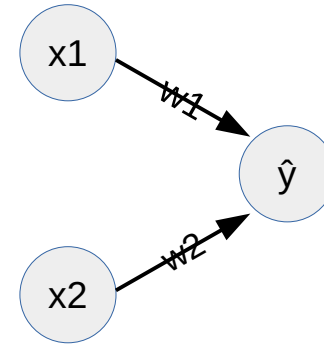
Backpropagation example, step by step:

<https://mattmazur.com/2015/03/17/a-step-by-step-backpropagation-example/>

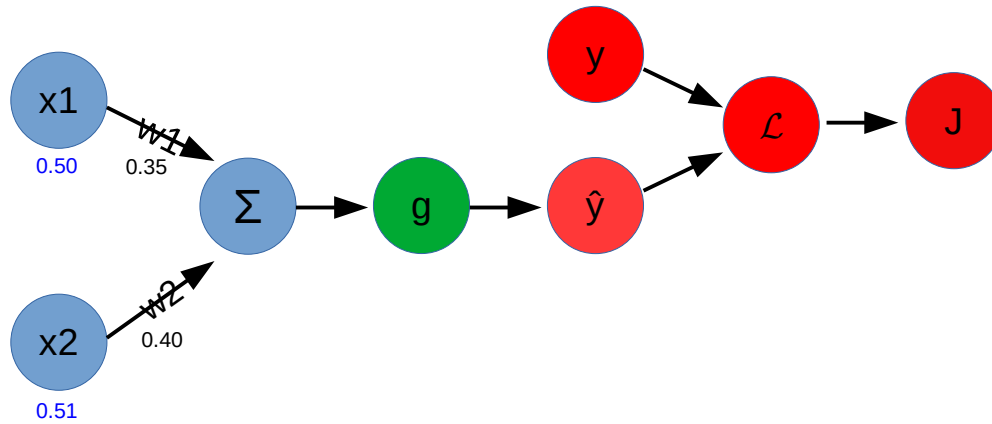
# Practical example

## Simple network:

- ▶ Two inputs  $[x_1, x_2]$
- ▶ Two weights  $[w_1, w_2]$
- ▶ No bias
- ▶ Activation function  $g()$
- ▶ One output  $\hat{y}$
- ▶ One label  $y$
- ▶ Loss function  $\mathcal{L}()$
- ▶ Weight-dependent error  $J(W)$



# 1. Forward pass



$$x_1 = 0.50$$

$$x_2 = 0.51$$

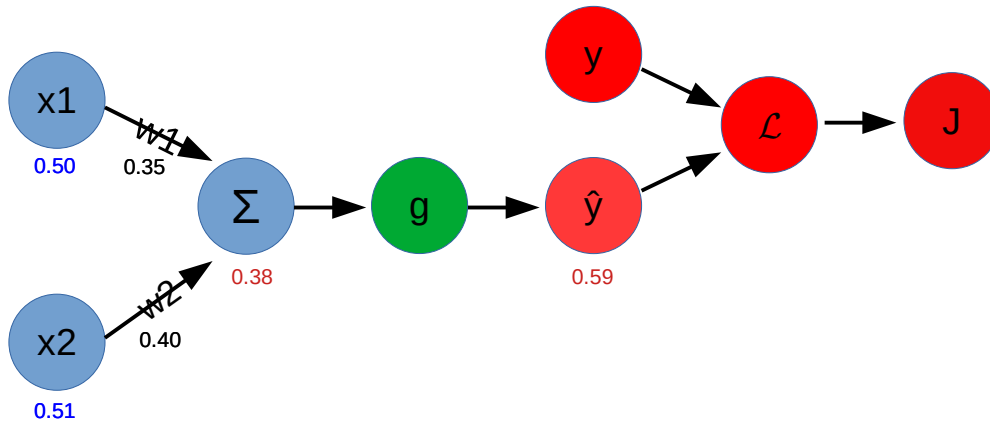
$$w_1 = 0.35$$

$$w_2 = 0.40$$

$$\Sigma = ?$$

$$\hat{y} = g(\Sigma) = 1/(1+e^{-\Sigma}) = ?$$

# 1. Forward pass



$$x_1 = 0.50$$

$$x_2 = 0.51$$

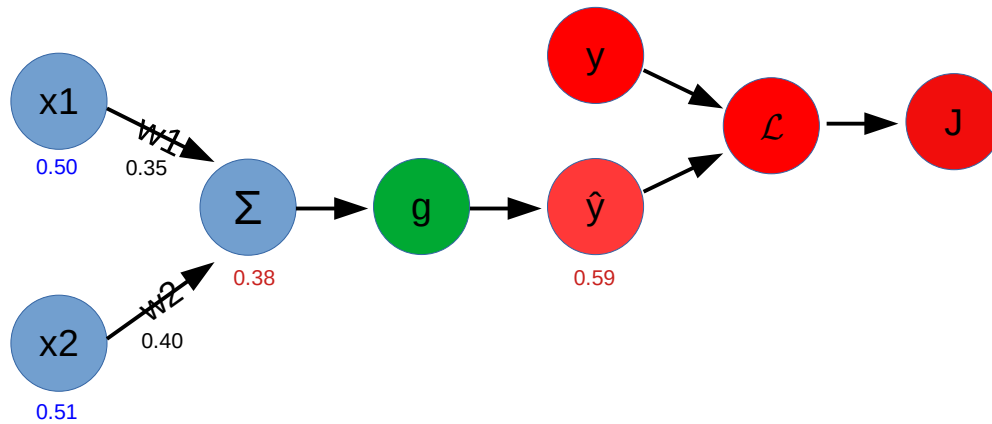
$$w_1 = 0.35$$

$$w_2 = 0.40$$

$$\Sigma = 0.35 * x_1 + 0.4 * w_2 = 0.38$$

$$\hat{y} = g(\Sigma) = 1/(1+e^{-\Sigma}) = 0.59$$

## 2. Calculate Loss



$$x_1 = 0.50$$

$$x_2 = 0.51$$

$$w_1 = 0.35$$

$$w_2 = 0.40$$

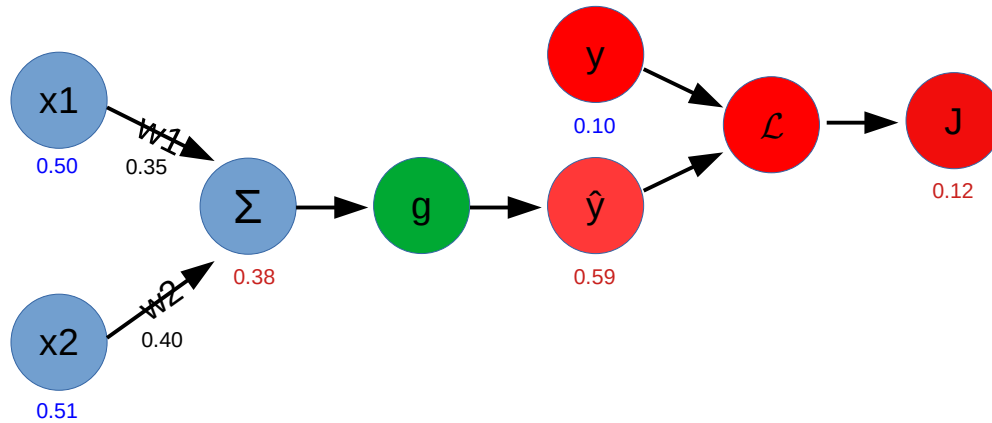
$$\Sigma = 0.35 * x_1 + 0.4 * w_2 = 0.38$$

$$\hat{y} = g(\Sigma) = 1/(1+e^{-\Sigma}) = 0.59$$

$$Y = 0.10$$

$$J(W) = \mathcal{L}(y, \hat{y}) = \frac{1}{2} * (y - \hat{y})^2 = ?$$

## 2. Calculate Loss



$$x_1 = 0.50$$

$$x_2 = 0.51$$

$$w_1 = 0.35$$

$$w_2 = 0.40$$

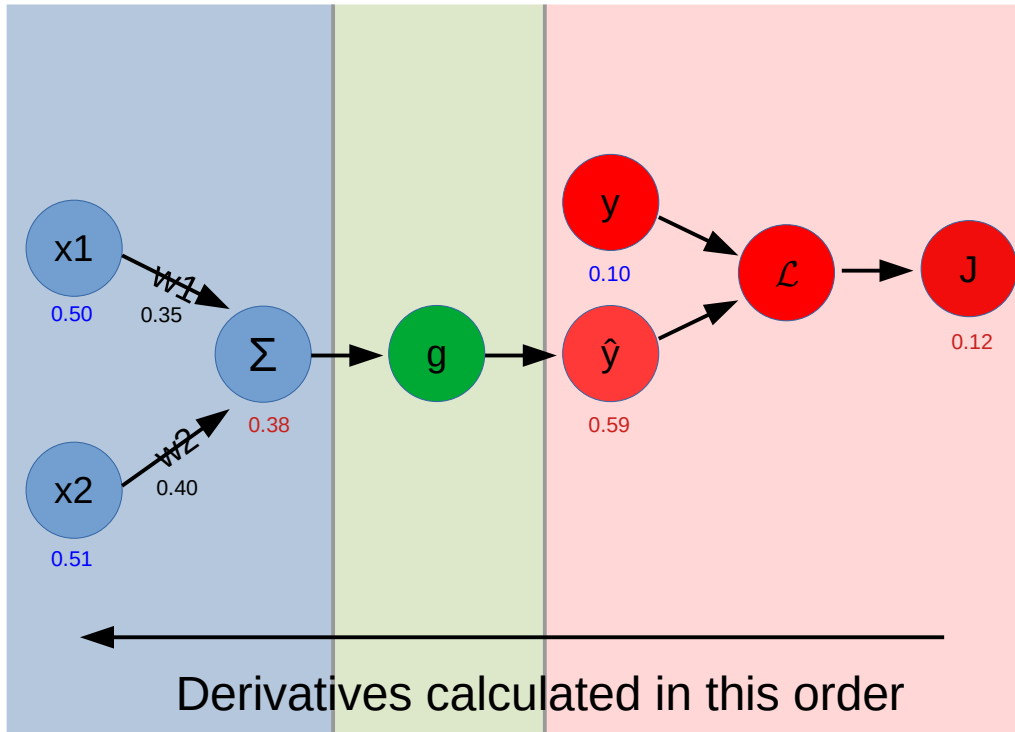
$$\Sigma = 0.35 * x_1 + 0.4 * w_2 = 0.38$$

$$\hat{y} = g(\Sigma) = 1/(1+e^{-\Sigma}) = 0.59$$

$$Y = 0.10$$

$$J(W) = \mathcal{L}(y, \hat{y}) = \frac{1}{2} * (y - \hat{y})^2 = 0.12$$

### 3. Backpropagate the error



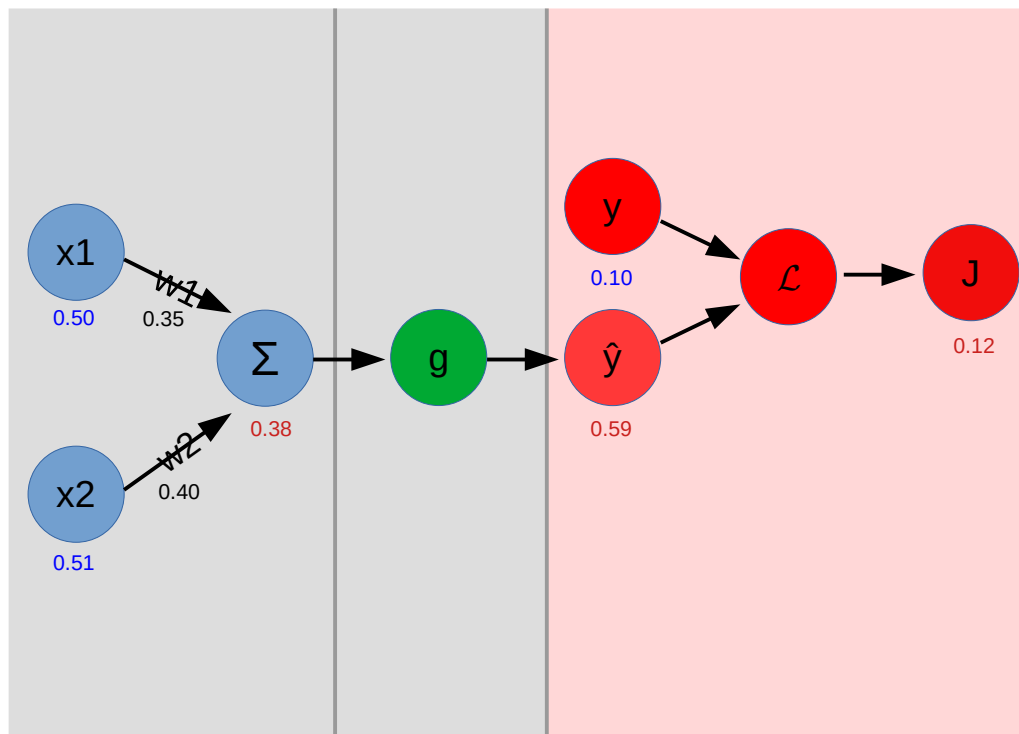
$$J(W) = \mathcal{L}(g(X * W))$$

$$\partial J(W) / \partial w_1 = \partial J(W) / \partial \hat{y} \quad (\text{loss})$$

$$* \partial \hat{y} / \partial \Sigma \quad (\text{activation})$$

$$* \partial \Sigma / \partial w_1 \quad (\text{weight})$$

### 3. Backpropagate the error: loss



$$J(W) = \mathcal{L}(g(X * W))$$

$$\partial J(W) / \partial w1 = \partial J(W) / \partial \hat{y} \quad (\text{loss})$$

$$* \partial \hat{y} / \Sigma \quad (\text{activation})$$

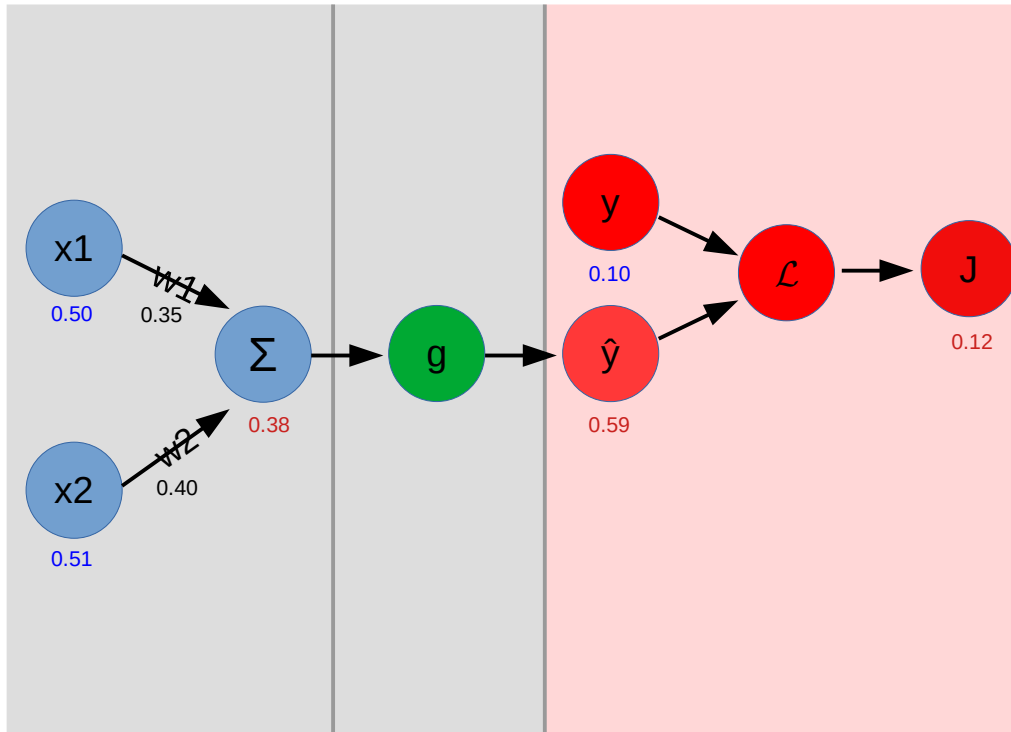
$$* \partial \Sigma / \partial w1 \quad (\text{weight})$$

$$J(W) = \mathcal{L}(y, \hat{y}) = \frac{1}{2} * (y - \hat{y})^2$$

$$\partial J(W) / \partial \hat{y} = \partial \mathcal{L}(y, \hat{y}) / \partial \hat{y} = ?$$



### 3. Backpropagate the error: loss



$$J(W) = \mathcal{L}(g(X * W))$$

$$\partial J(W) / \partial w_1 = \partial J(W) / \partial \hat{y} \quad (\text{loss})$$

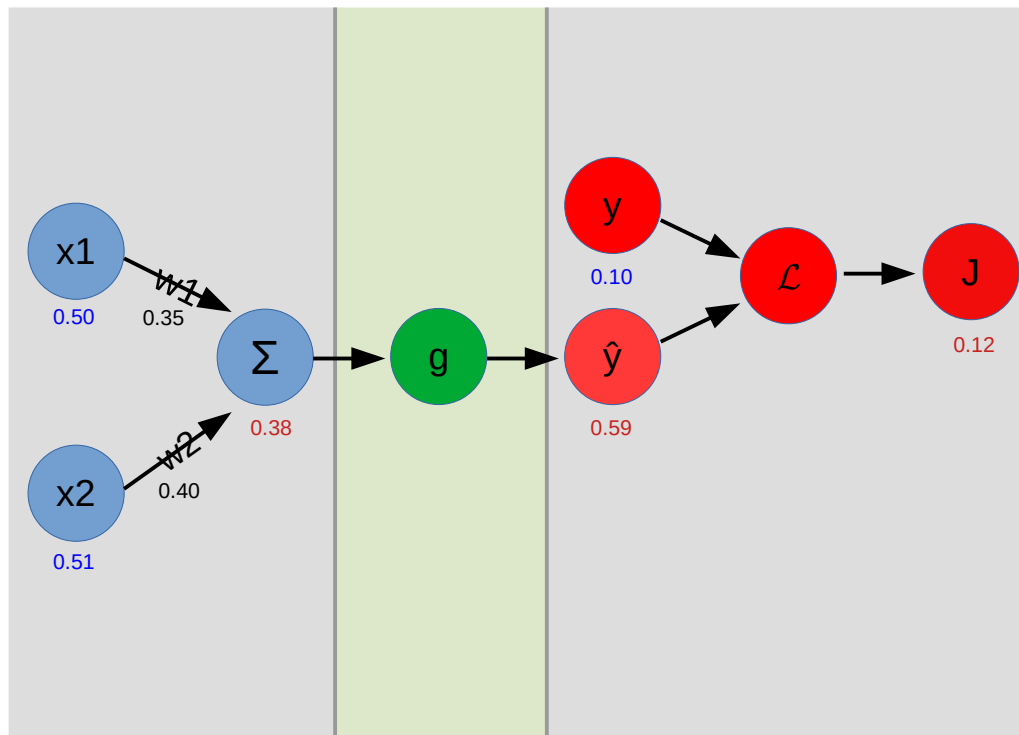
$$* \partial \hat{y} / \partial \Sigma \quad (\text{activation})$$

$$* \partial \Sigma / \partial w_1 \quad (\text{weight})$$

$$J(W) = \mathcal{L}(y, \hat{y}) = \frac{1}{2} * (y - \hat{y})^2$$

$$\begin{aligned} \partial J(W) / \partial \hat{y} &= 2 * \frac{1}{2} * (y - \hat{y}) * -1 \\ &= -y + \hat{y} = -0.1 + 0.59 = 0.49 \end{aligned}$$

### 3. Backpropagate the error: activation



$$J(W) = \mathcal{L}(g(X * W))$$

$$\partial J(W) / \partial w1 = 0.49 \quad (\text{loss})$$

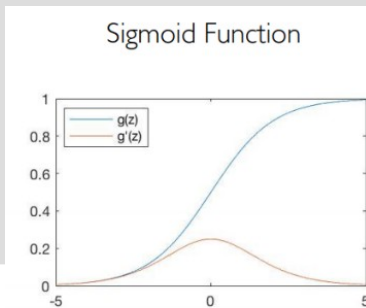
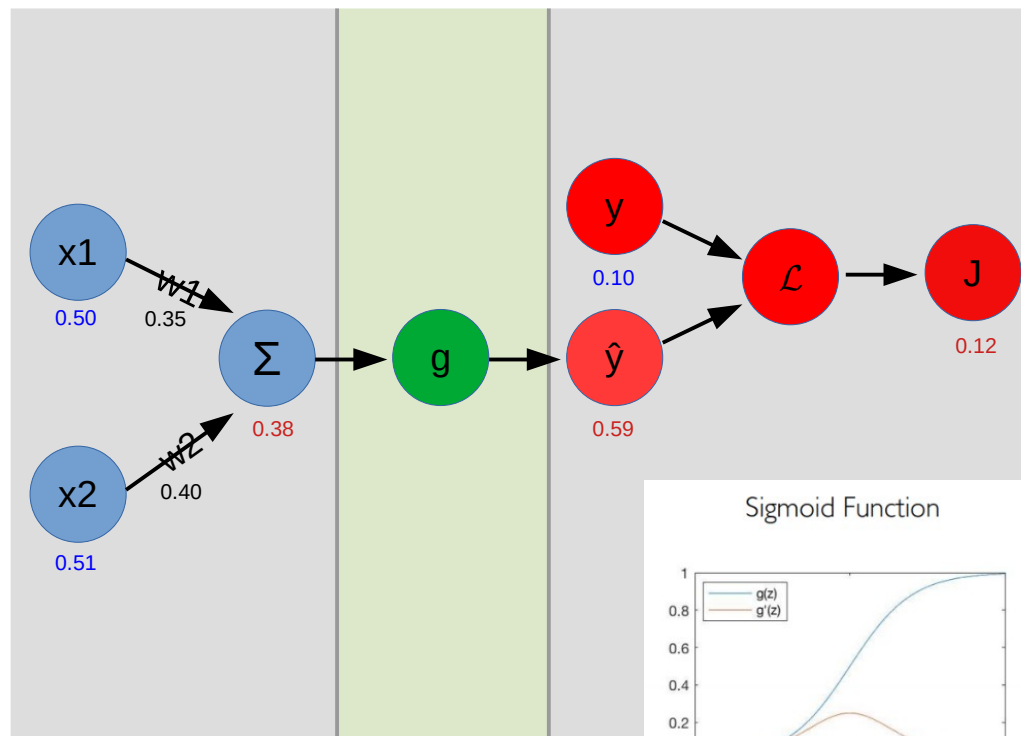
$$* \partial \hat{y} / \partial \Sigma \quad (\text{activation})$$

$$* \partial \Sigma / \partial w1 \quad (\text{weight})$$

$$\hat{y} = g(\Sigma) = 1/(1+e^{-\Sigma})$$

$$\partial \hat{y} / \partial \Sigma = ?$$

### 3. Backpropagate the error: activation



$$g(z) = \frac{1}{1 + e^{-z}}$$

$$g'(z) = g(z)(1 - g(z))$$

$$J(W) = \mathcal{L}(g(X * W))$$

$$\partial J(W) / \partial w_1 = 0.49 \quad (\text{loss})$$

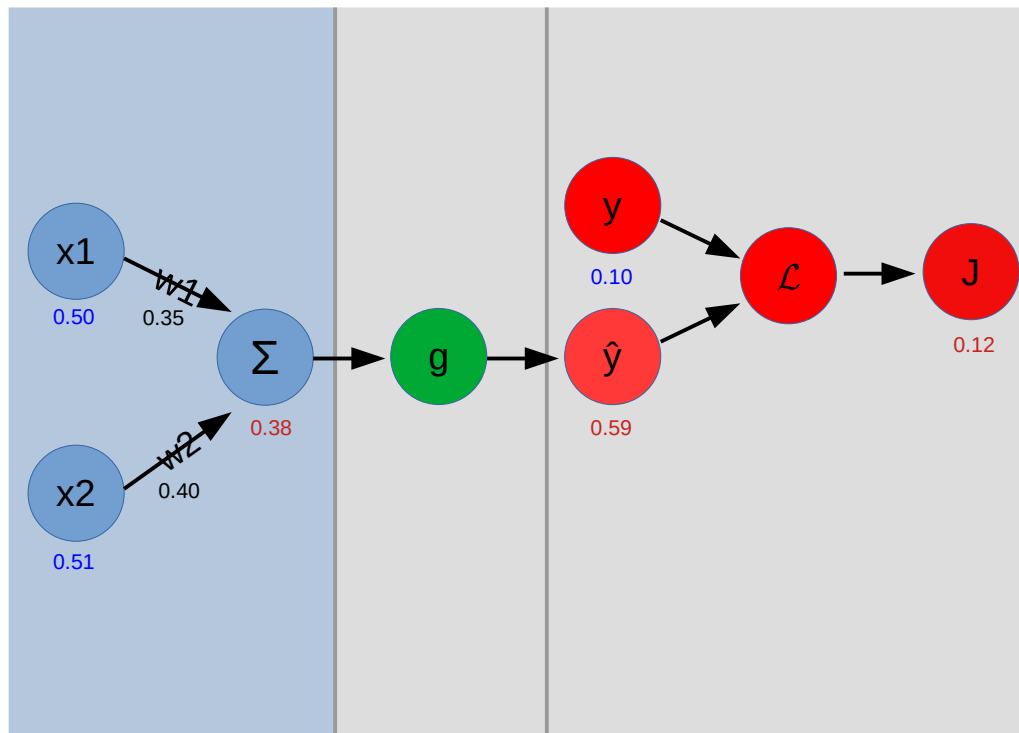
$$* \partial \hat{y} / \partial \Sigma \quad (\text{activation})$$

$$* \partial \Sigma / \partial w_1 \quad (\text{weight})$$

$$\hat{y} = g(\Sigma) = 1/(1 + e^{-\Sigma})$$

$$\begin{aligned} \partial \hat{y} / \partial \Sigma &= 1/(1 + e^{-\Sigma}) * (1 - 1/(1 + e^{-\Sigma})) \\ &= 0.24 \end{aligned}$$

### 3. Backpropagate the error: weight



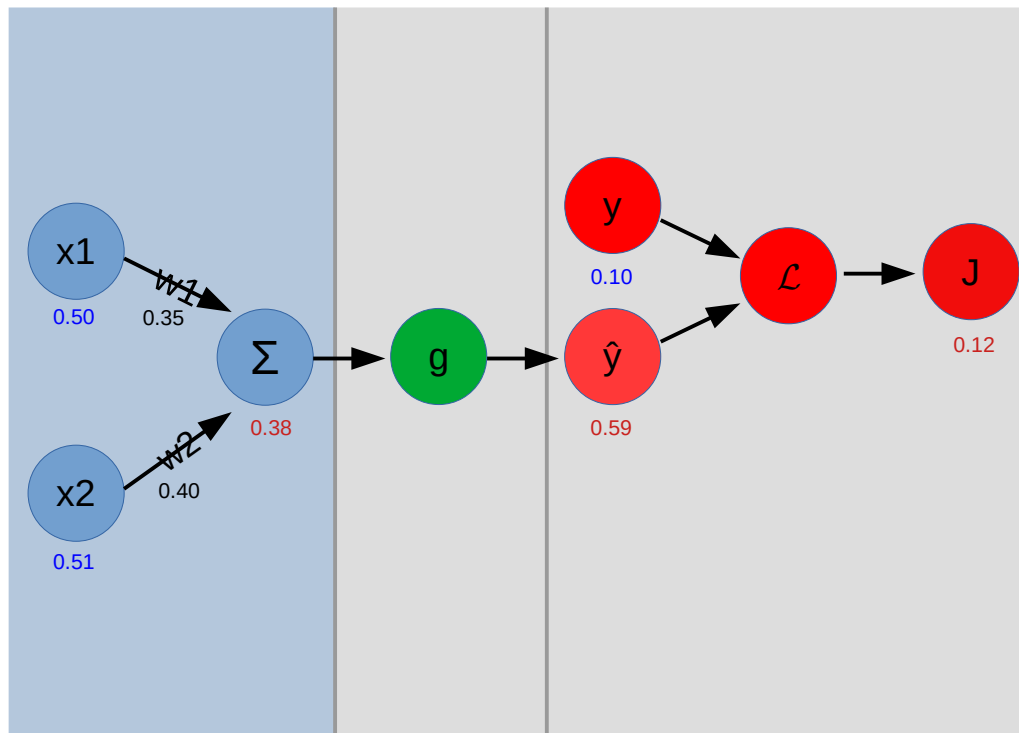
$$J(W) = \mathcal{L}(g(X * W))$$

$$\begin{aligned} \partial J(W) / \partial w1 &= 0.49 && \text{(loss)} \\ &* 0.24 && \text{(activation)} \\ &* \partial \Sigma / \partial w1 && \text{(weight)} \end{aligned}$$

$$\Sigma = X * W = x1 * w1 + x2 * w2$$

$$\partial \Sigma / \partial w1 = ?$$

### 3. Backpropagate the error: weight



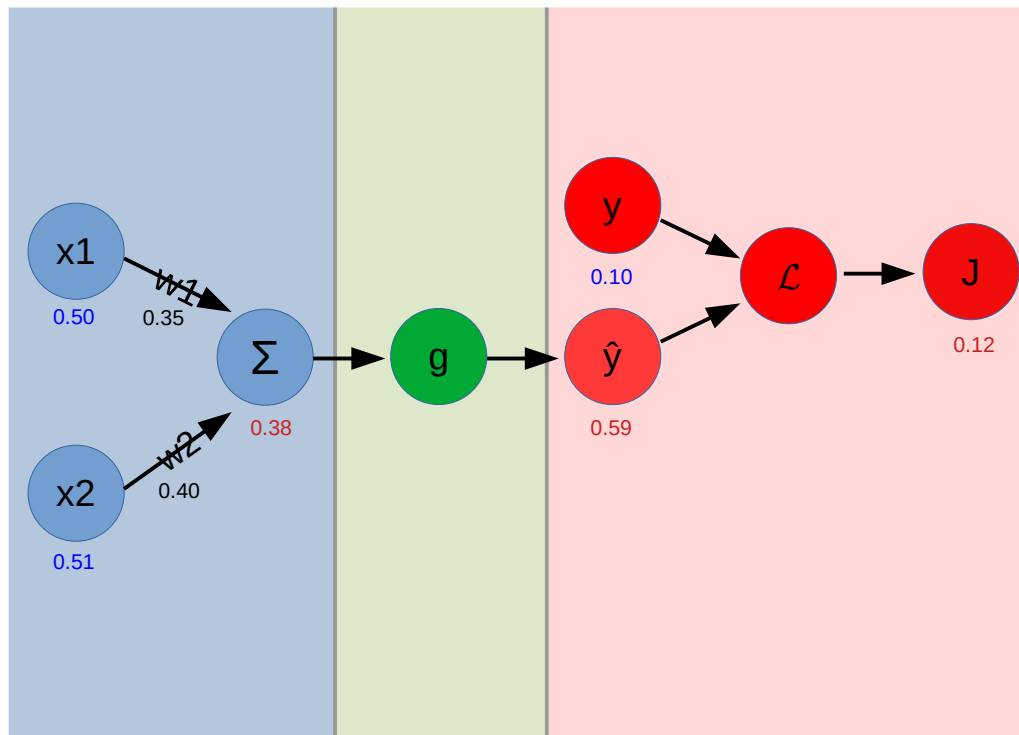
$$J(W) = \mathcal{L}(g(X * W))$$

$$\begin{aligned} \frac{\partial J(W)}{\partial w_1} &= 0.49 && \text{(loss)} \\ &\quad * 0.24 && \text{(activation)} \\ &\quad * \frac{\partial \Sigma}{\partial w_1} && \text{(weight)} \end{aligned}$$

$$\Sigma = X * W = x_1 * w_1 + x_2 * w_2$$

$$\frac{\partial \Sigma}{\partial w_1} = x_1 + 0 = 0.5$$

## 4. Weight update



$$J(W) = \mathcal{L}(g(X * W))$$

$$\partial J(W) / \partial w_1 = 0.49$$

$$* 0.24$$

$$* 0.51$$

$$= 0.06$$

$$w_1' = ?$$

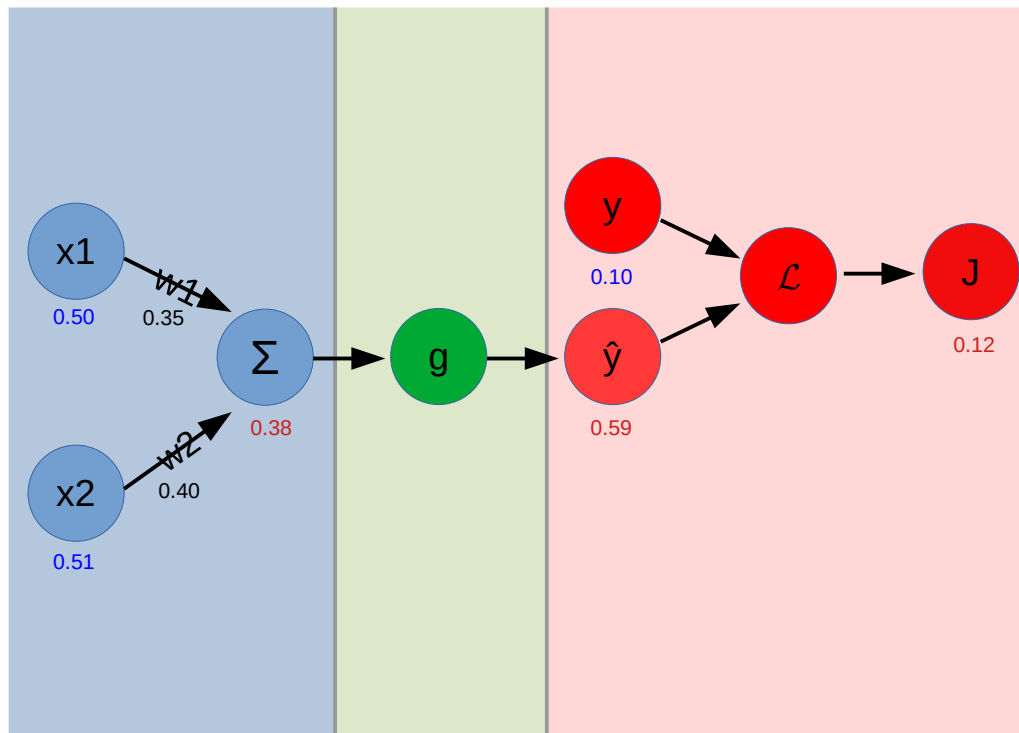
(loss)

(activation)

(weight)

(gradient)

## 4. Weight update



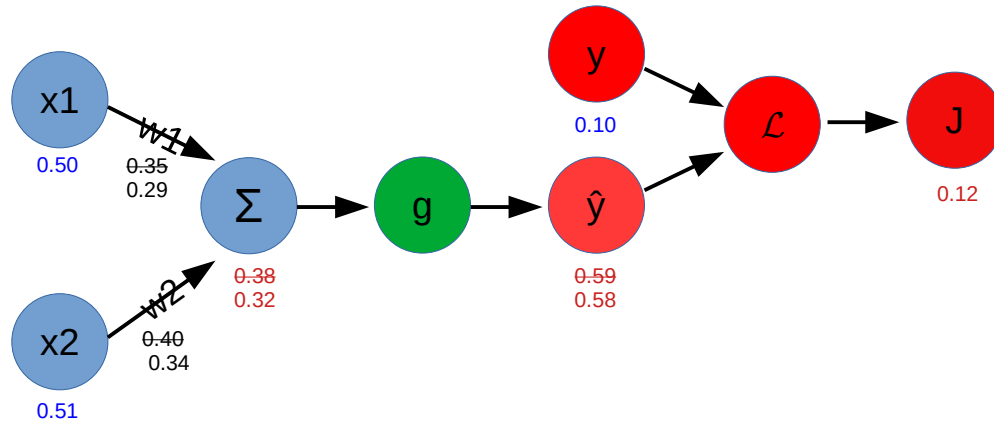
$$J(W) = \mathcal{L}(g(X * W))$$

$$\begin{aligned} \frac{\partial J(W)}{\partial w_1} &= 0.49 && \text{(loss)} \\ &\quad * 0.24 && \text{(activation)} \\ &\quad * 0.51 && \text{(weight)} \\ &= 0.06 && \text{(gradient)} \end{aligned}$$

$$w_1' = w_1 - \eta * 0.06 = 0.35 - 0.06 = 0.29$$

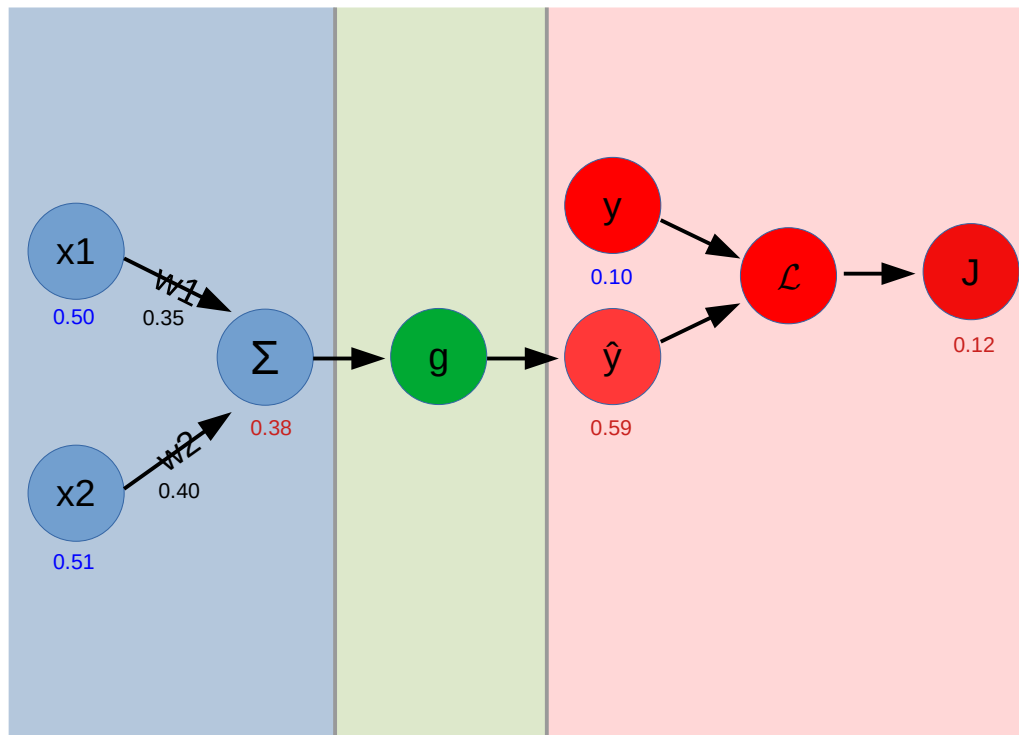
# Exercise

- ▶ Can you calculate the weight update for  $w_2$ ? How many new gradients do you need to calculate?
- ▶ What is the new predicted output? Has the error gone down?
- ▶ What if I had another layer before this one?





## 4. Weight update



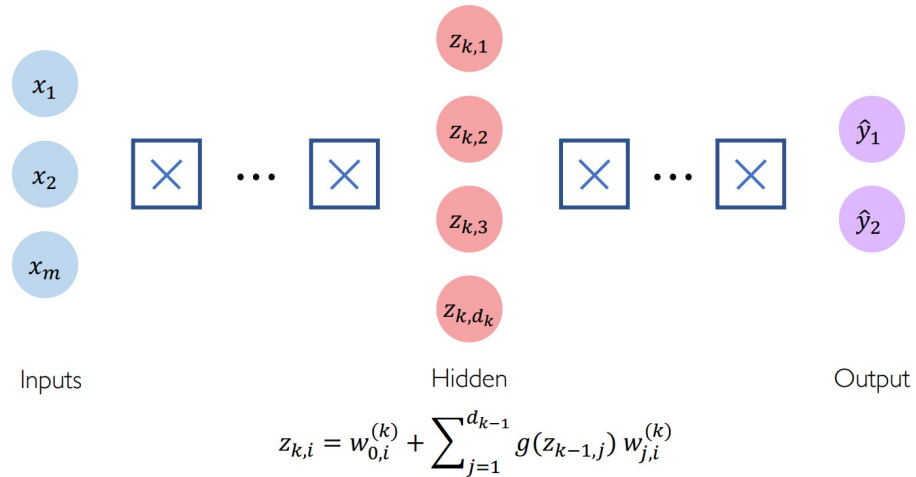
$$J(W) = \mathcal{L}(g(X * W))$$

$$\begin{aligned} \partial J(W) / \partial w_1 &= 0.49 && \text{(loss)} \\ &\quad * 0.24 && \text{(activation)} \\ &\quad * 0.51 && \text{(weight)} \\ &= 0.06 && \text{(gradient)} \end{aligned}$$

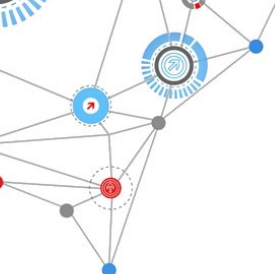
$$\begin{aligned} w_1' &= w_1 - \eta * 0.06 = 0.35 - 0.06 = 0.29 \\ w_2' &= w_2 - \eta * 0.06 = 0.4 - 0.06 = 0.34 \end{aligned}$$

# Gradient vanishing

- ▶ What happens if we backpropagate on a network with many ( $N > k > 1$ ) hidden layers?



$$\partial J / \partial W_1 = \partial J / \partial \hat{Y} * \partial \hat{Y} / \partial \Sigma_N * \partial \Sigma_N / \partial W_N * \dots * \partial Z_k / \partial \Sigma_k * \partial \Sigma_k / \partial W_k * \dots * \partial Z_1 / \partial \Sigma_1 * \partial \Sigma_1 / \partial W_1$$



# Gradient vanishing

- ▶ These are all “zero-point-somethings” multiplied by each other
- ▶ So the gradient becomes smaller by orders of magnitudes as we go back more and more layers until it's so small that the network is stuck

$$\partial E / \partial W_1 = \partial J / \partial \hat{Y} * \partial \hat{Y} / \partial \Sigma_N * \partial \Sigma_N / \partial W_N * \dots * \partial Z_k / \partial \Sigma_k * \partial \Sigma_k / \partial W_k * \dots * \partial Z_1 / \partial \Sigma_1 * \partial \Sigma_1 / \partial X_1 = O(10^{-N})$$

initial  $w_1 = 0.5$

optimal  $w_1 = -0.2$

5-layer gradient  $\sim 0.00001$

How many iterations do we need to get from 0.5 to -0.2?