

# *Neural Network 101*

Rafael Geurgas

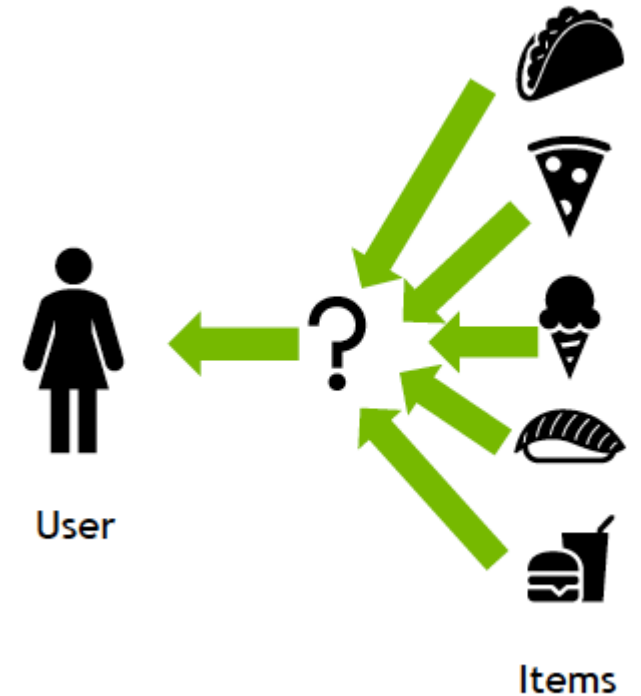
02/07/25



SCAN ME

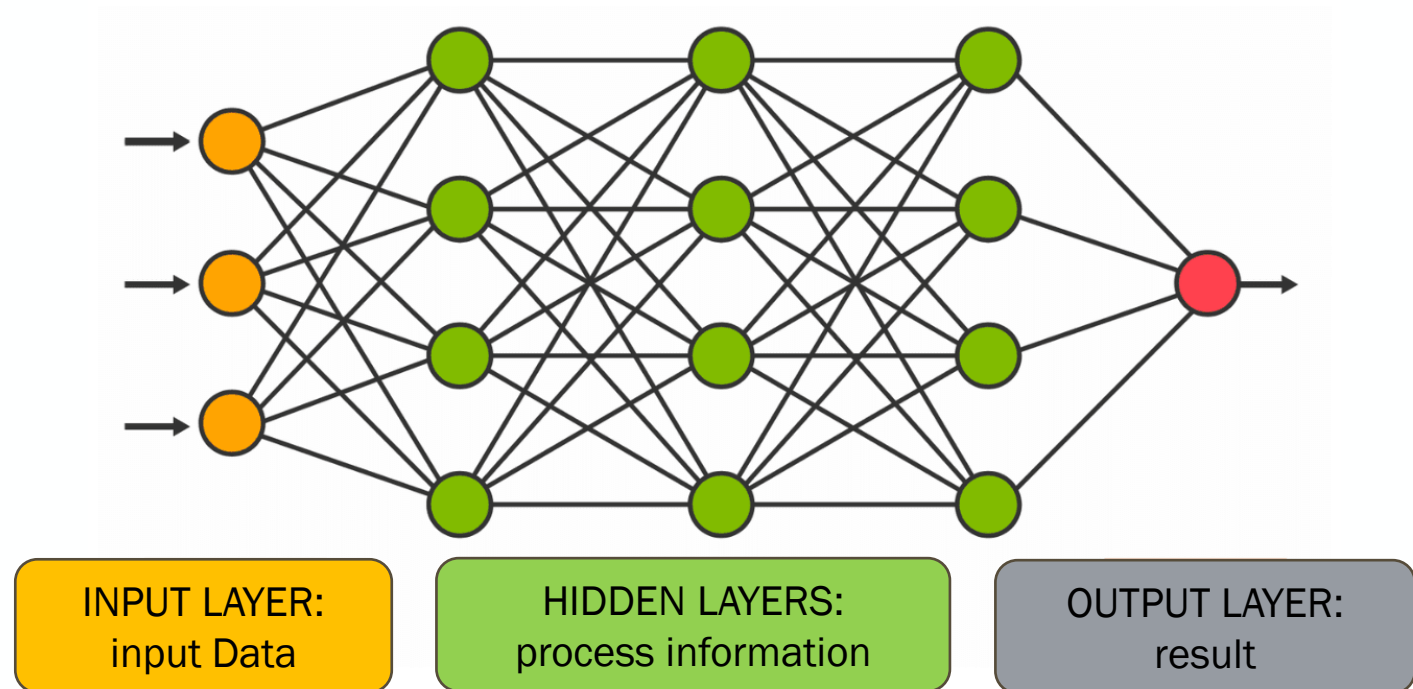
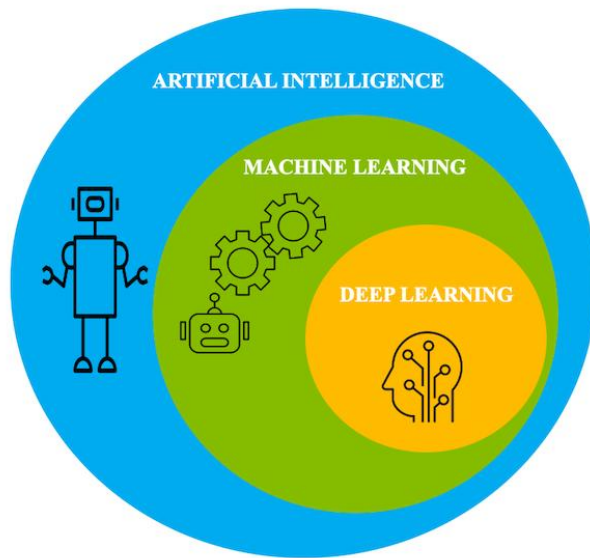
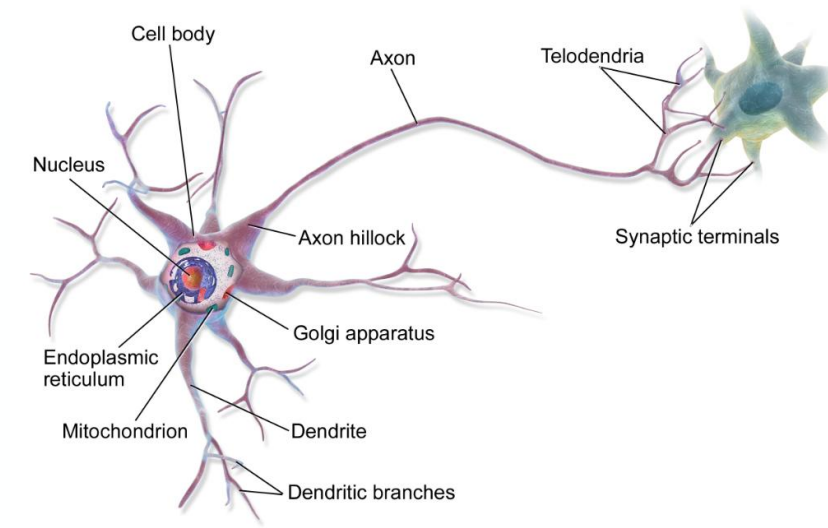
# Artificial Intelligence

- Artificial intelligence (AI): imitating human intelligence or behavioral patterns or any other living entity;
- Machine Learning (ML) is a branch of AI that focuses on the development of algorithms and statistical models that enable computers to perform tasks without being explicitly programmed for those tasks.
- ML is commonly used in problems where the answer is easily obtainable, but explaining it explicitly is not.
  - Pattern Recognition;
  - Recommendation Systems.

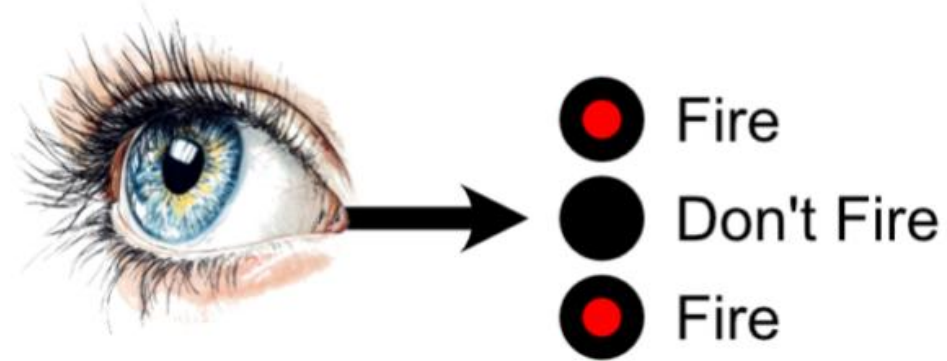


# Deep Learning

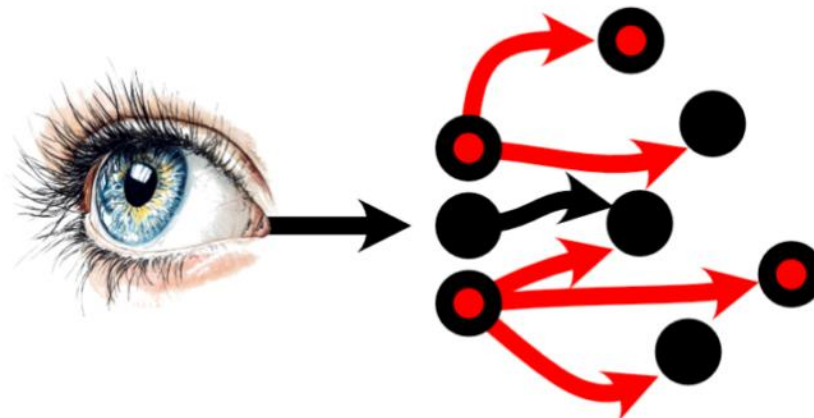
Deep learning is a subset of AI and machine learning (ML) that aims to mimic the human brain's structure and function through **artificial neural networks**.



# *How the Brain Works*

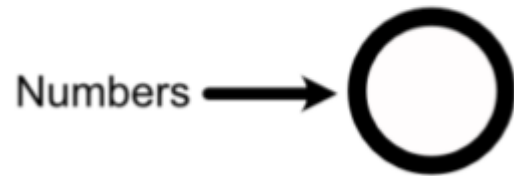


Imagine the signal from the eye directly feeds into three neurons, and two decide to fire.

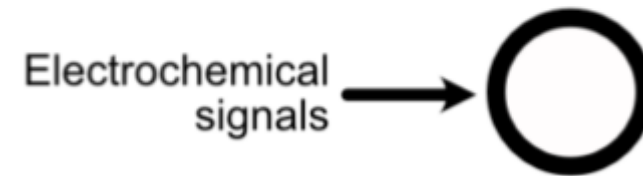


# *Similarity between Neural Networks and the Brain*

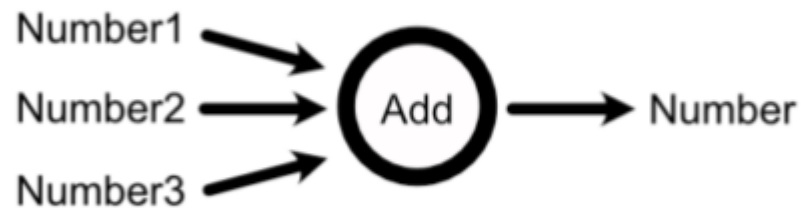
A Perceptron



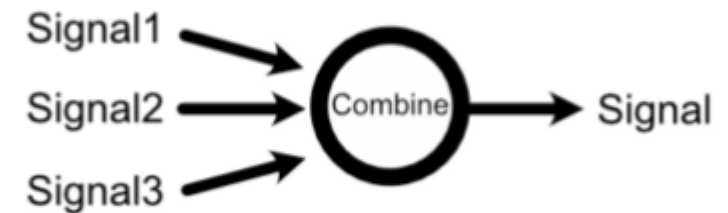
A Neuron



A Perceptron



A Neuron



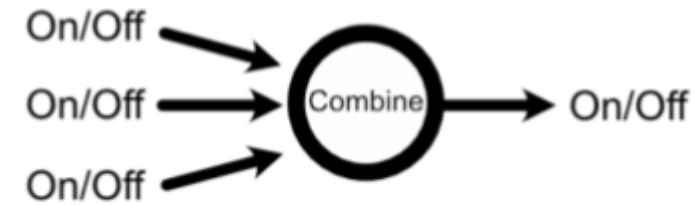
Perceptrons output numbers, while neurons output electrochemical signals.

# *Similarity between Neural Networks and the Brain*

A Perceptron



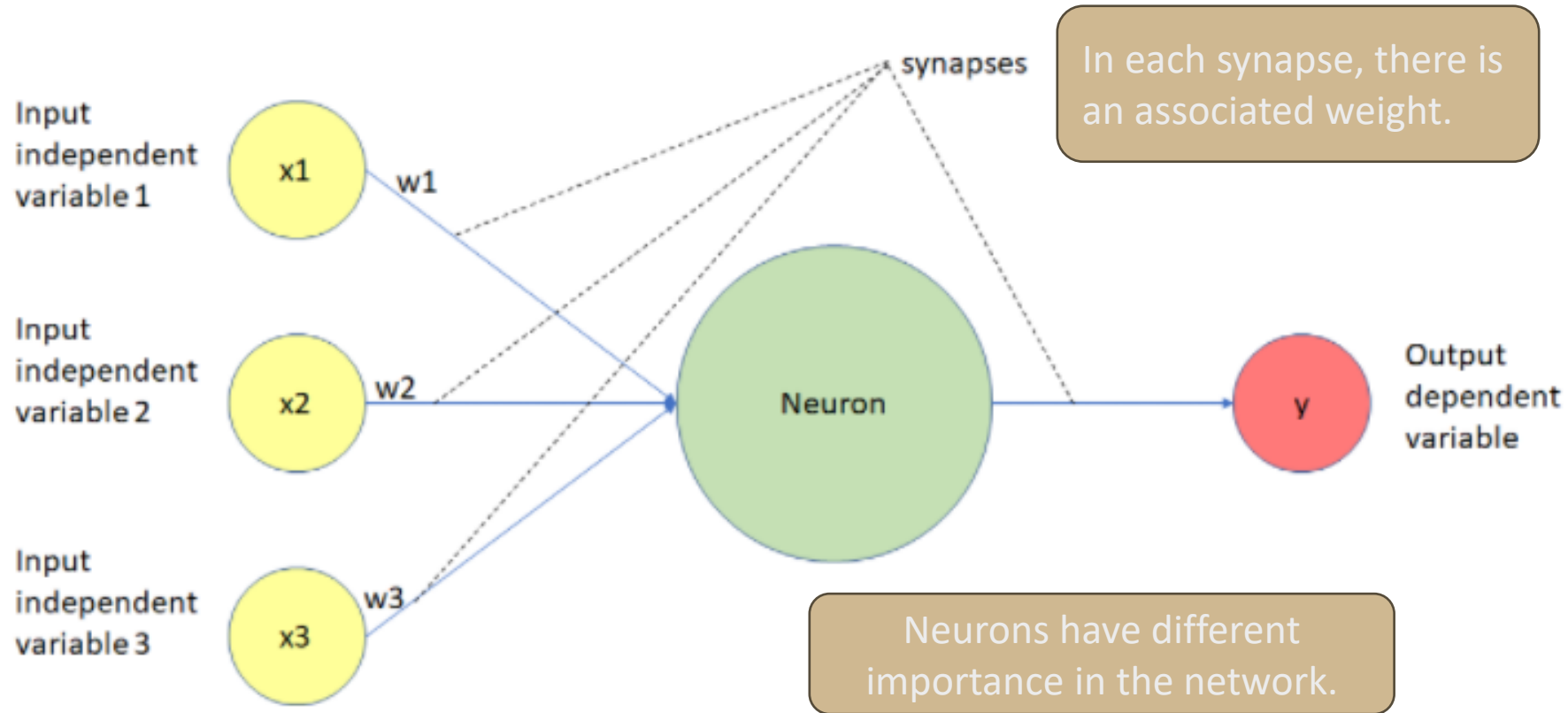
A Neuron



Perceptrons output a continuous range of numbers, while Neurons either fire or they don't.

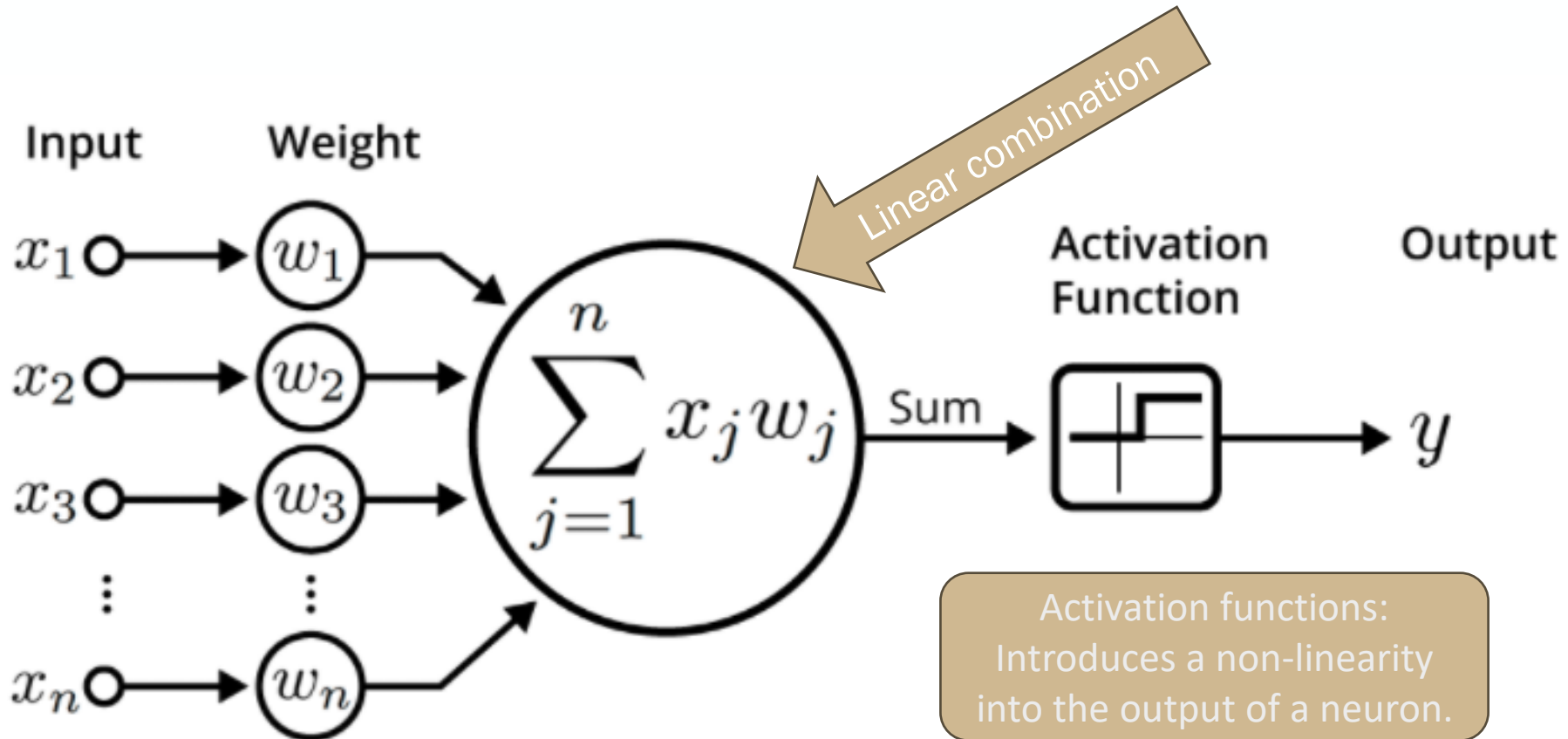
# Neural Network: Neuron (Perceptron)

Node where information and computational calculations are passed.



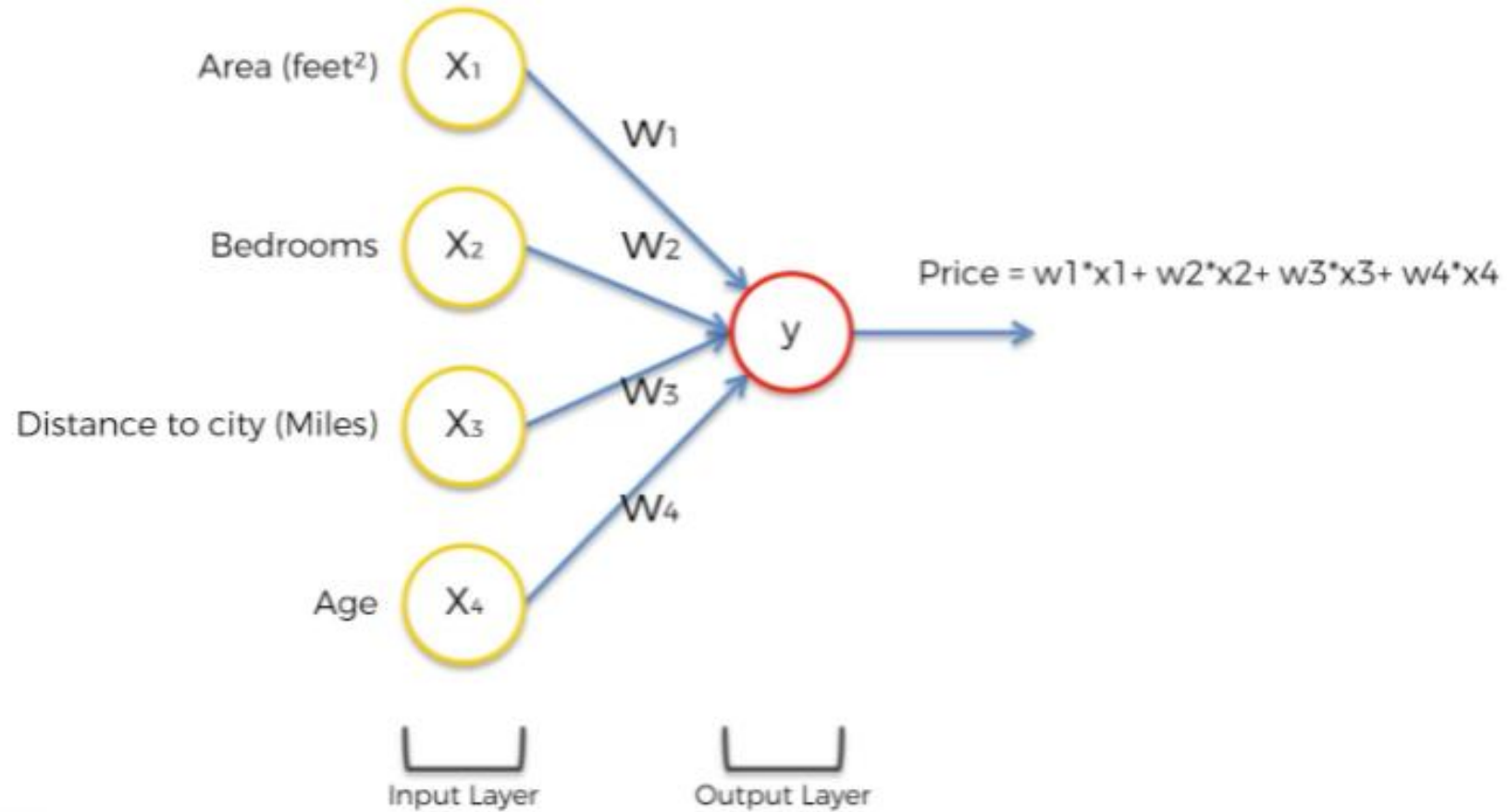
# Weights

The progressive adjustment of weights is the primary way in which a neural network learns.

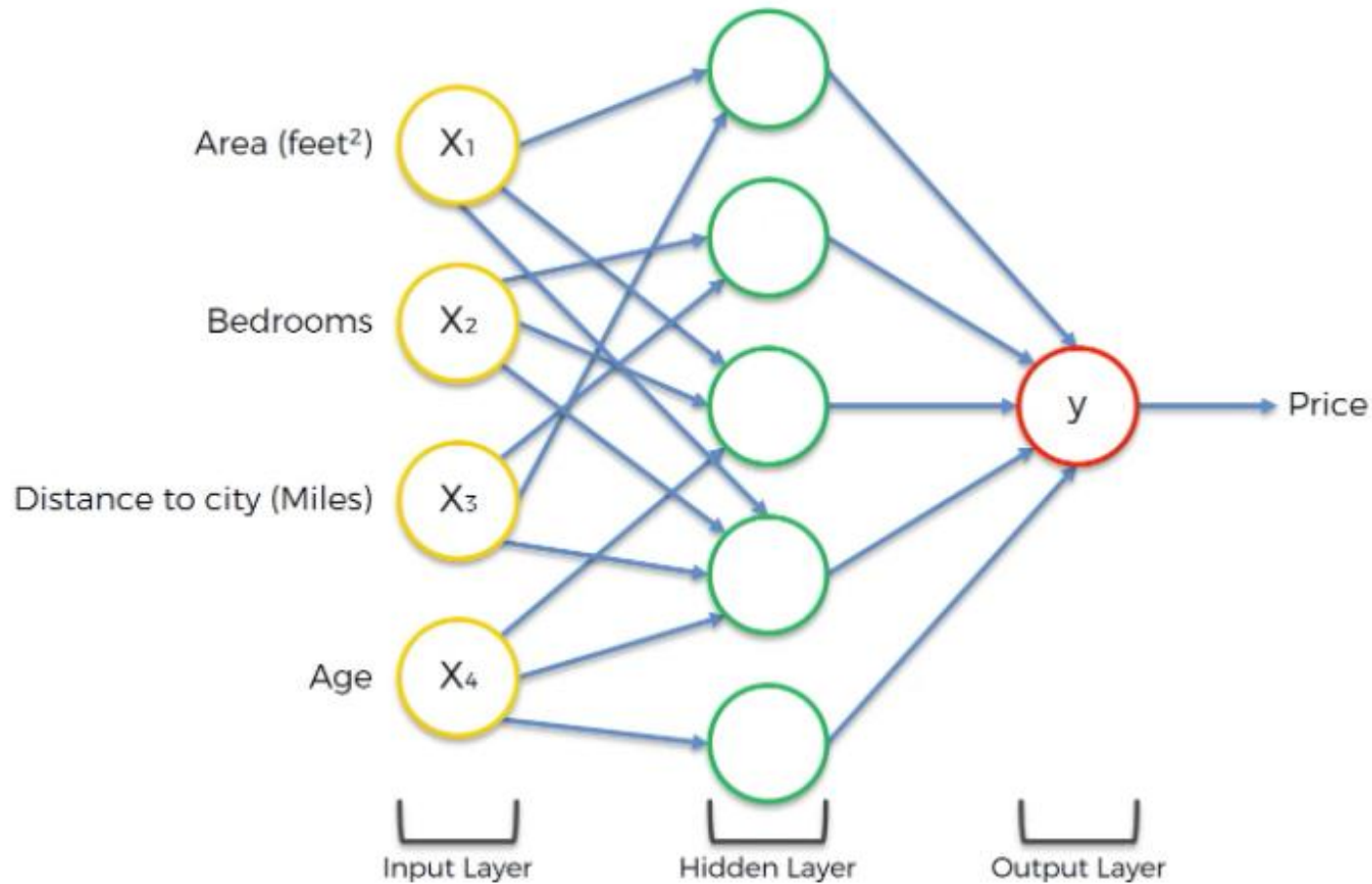




# Simple Example



# Simple Example



- Adding a Hidden layer adds more parameters (weights) to the network.
- From 4 to 16.
- Adding an activation function can give more freedom to the output

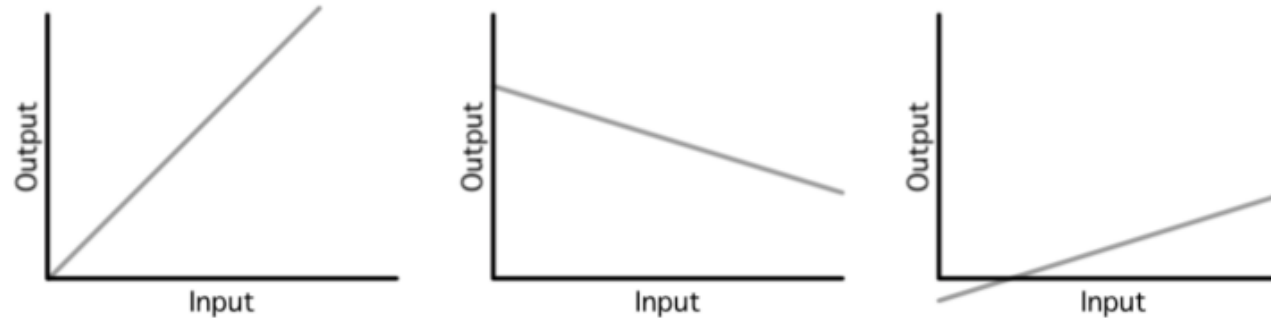
# Bias



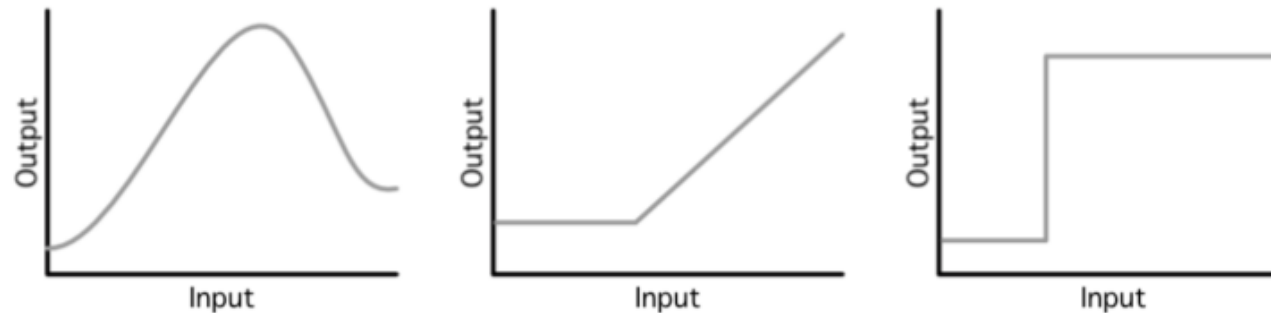
# Activation Function

- This is the biggest difference between Neural Network and Linear Regression

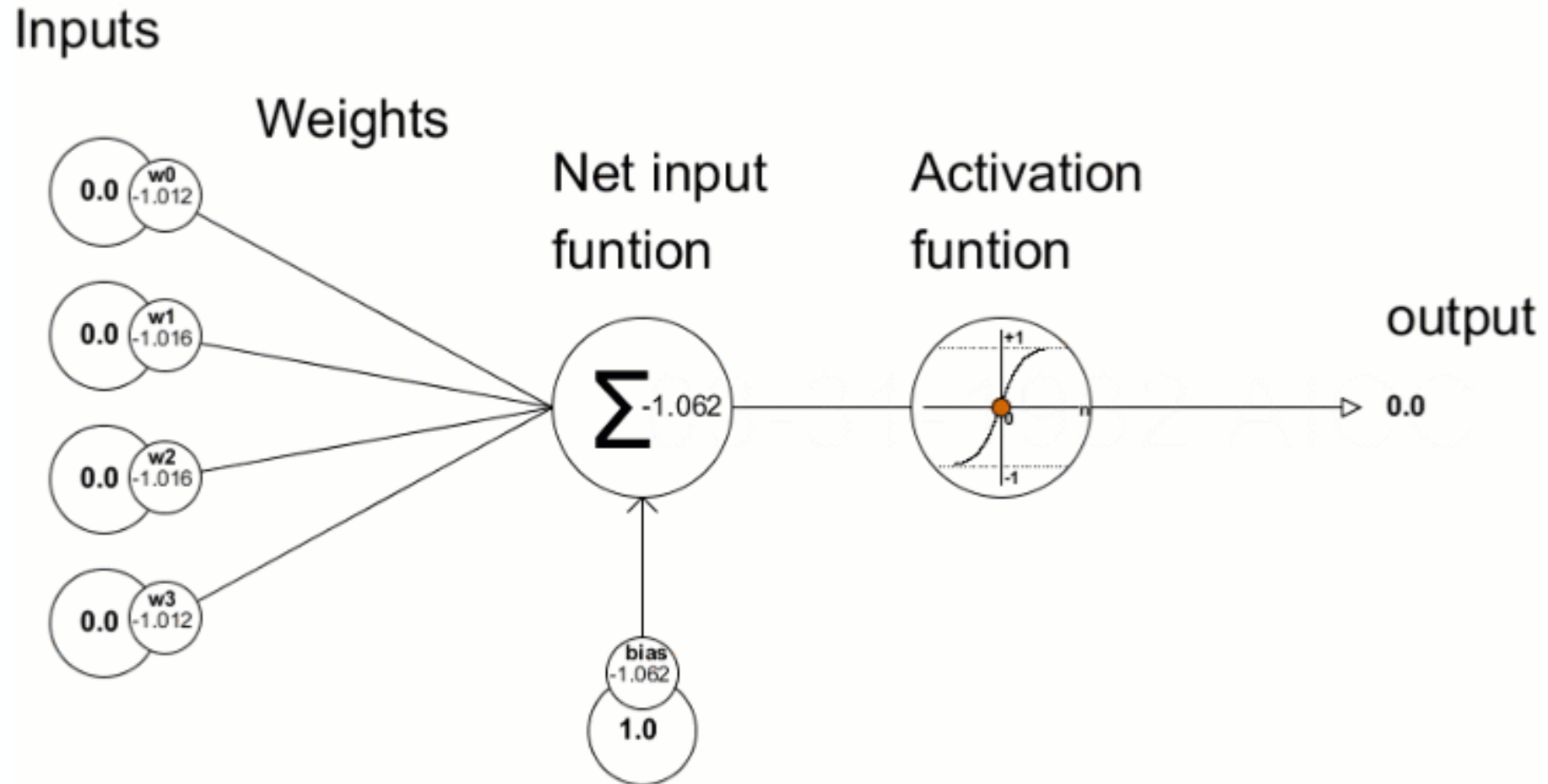
## Linear Functions



## Non-Linear Functions



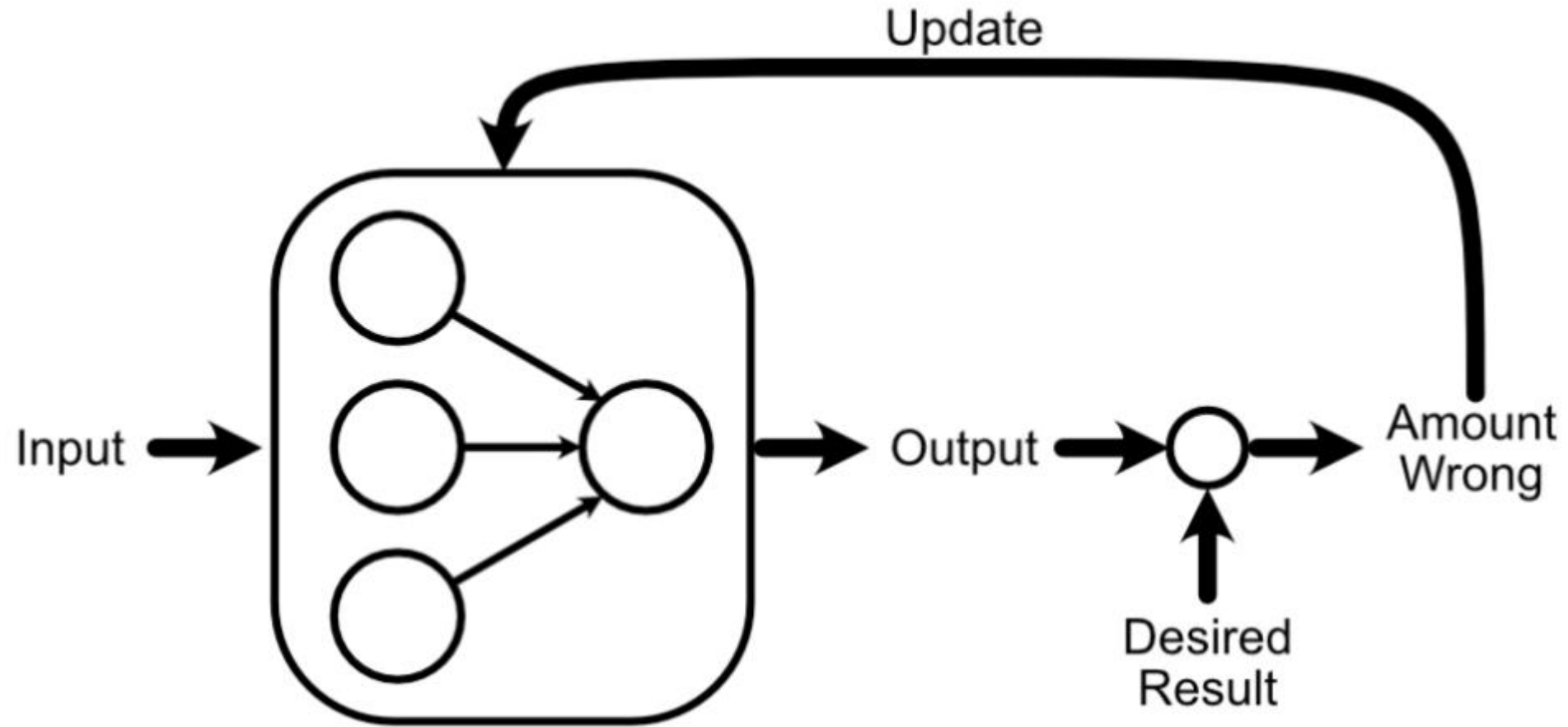
# Neuron

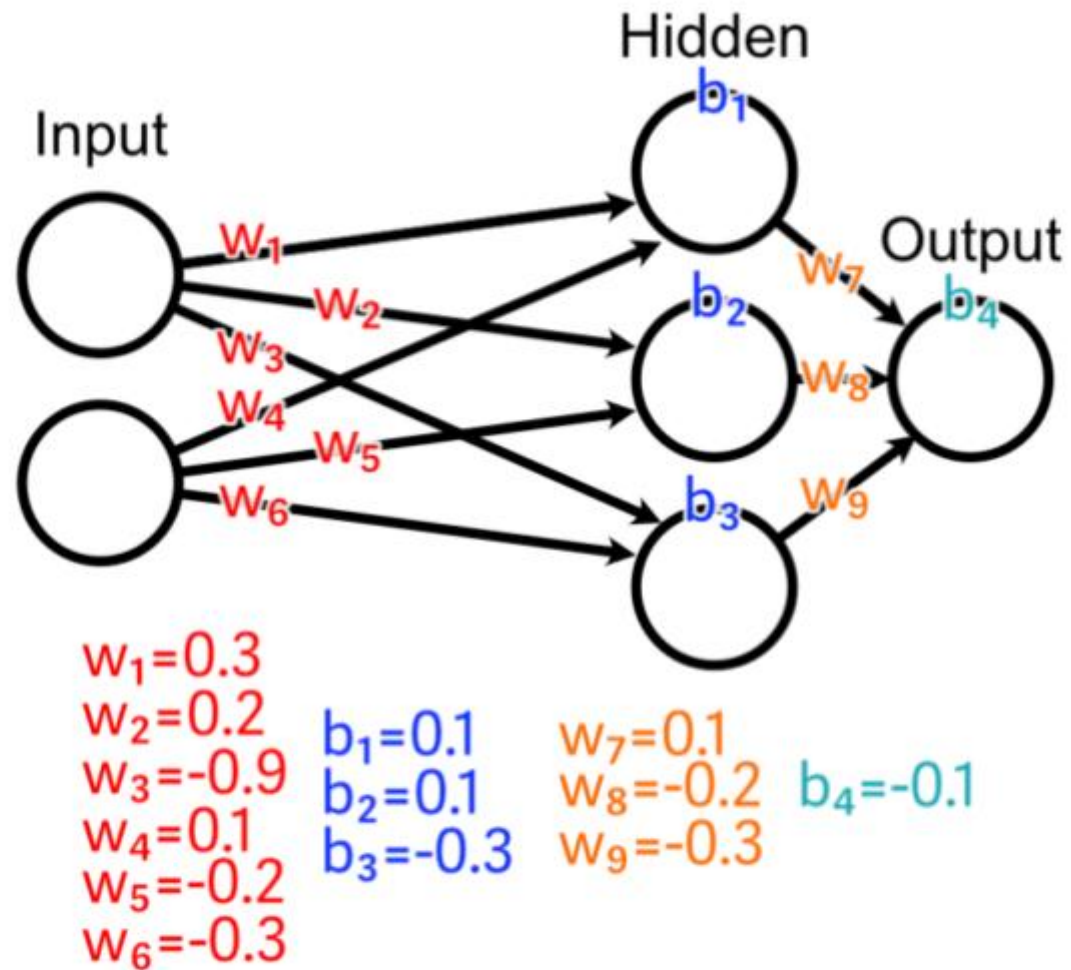


# *Classification*

- Unsupervised: There is no information about classification in the database.
  - Clustering;
  - Recommendation Algorithms.
- Supervised: There is information about classification in the database.
  - Classification;
  - Prediction.
- Semi-supervised: There is information about classification in a fraction of data in the database.

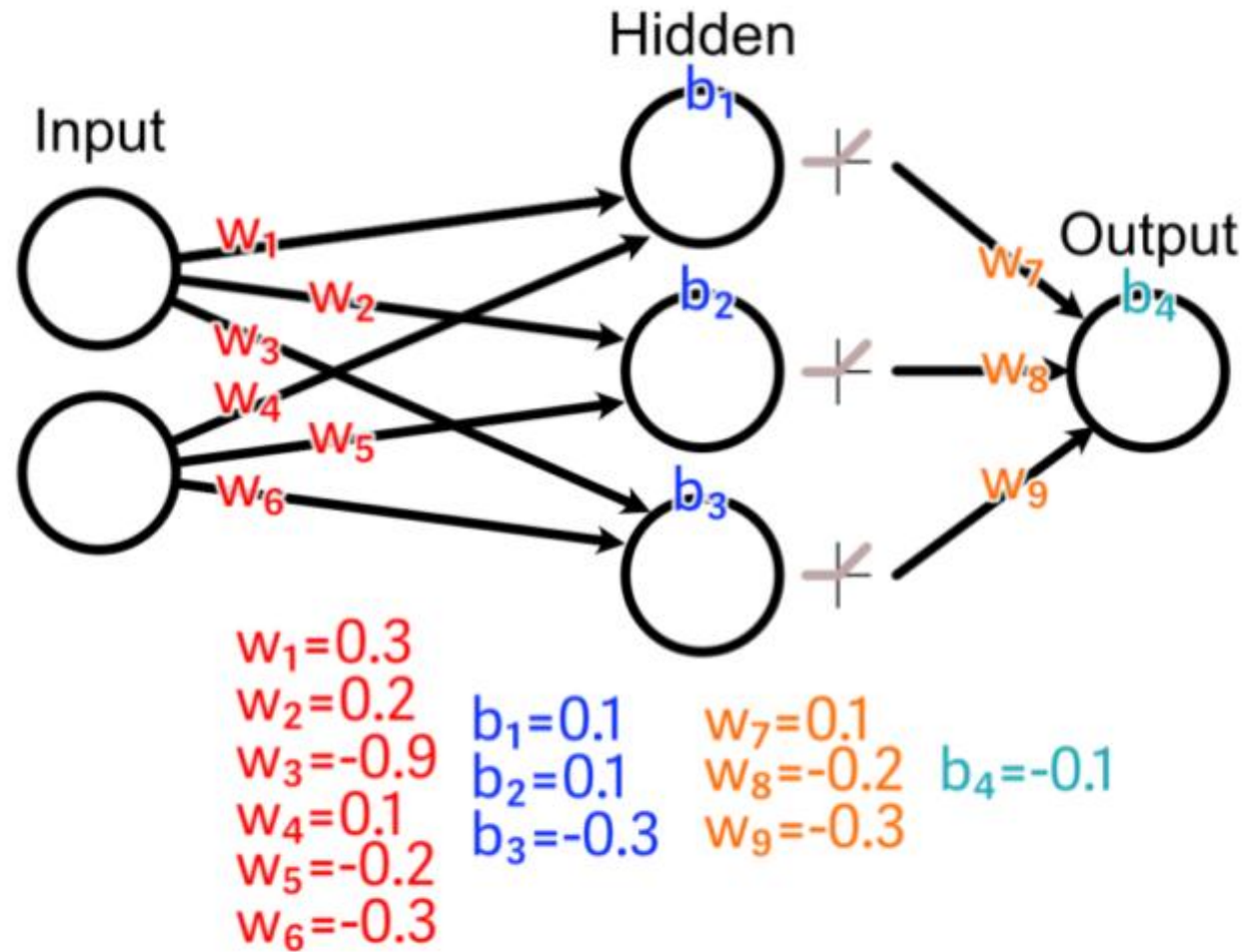
# *Training Process*





The neural network, with random weights and biases defined.

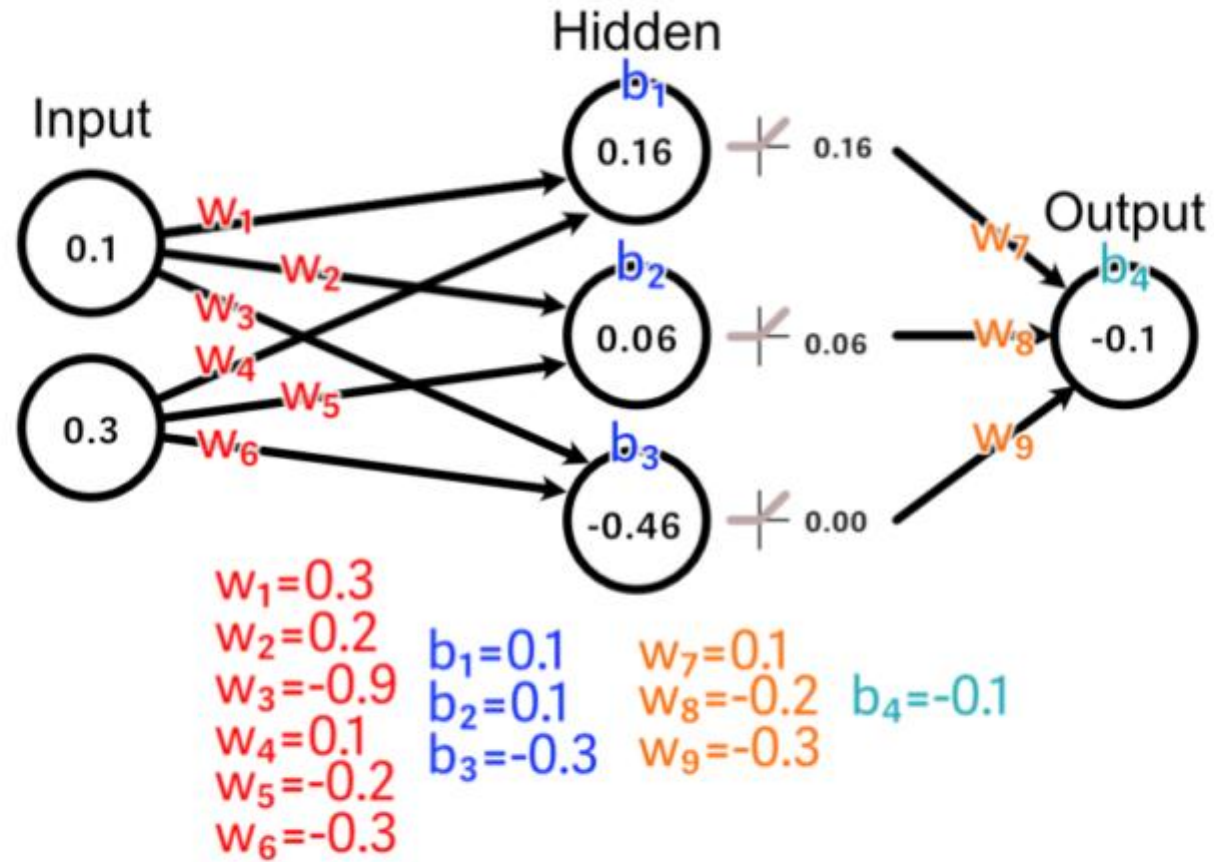




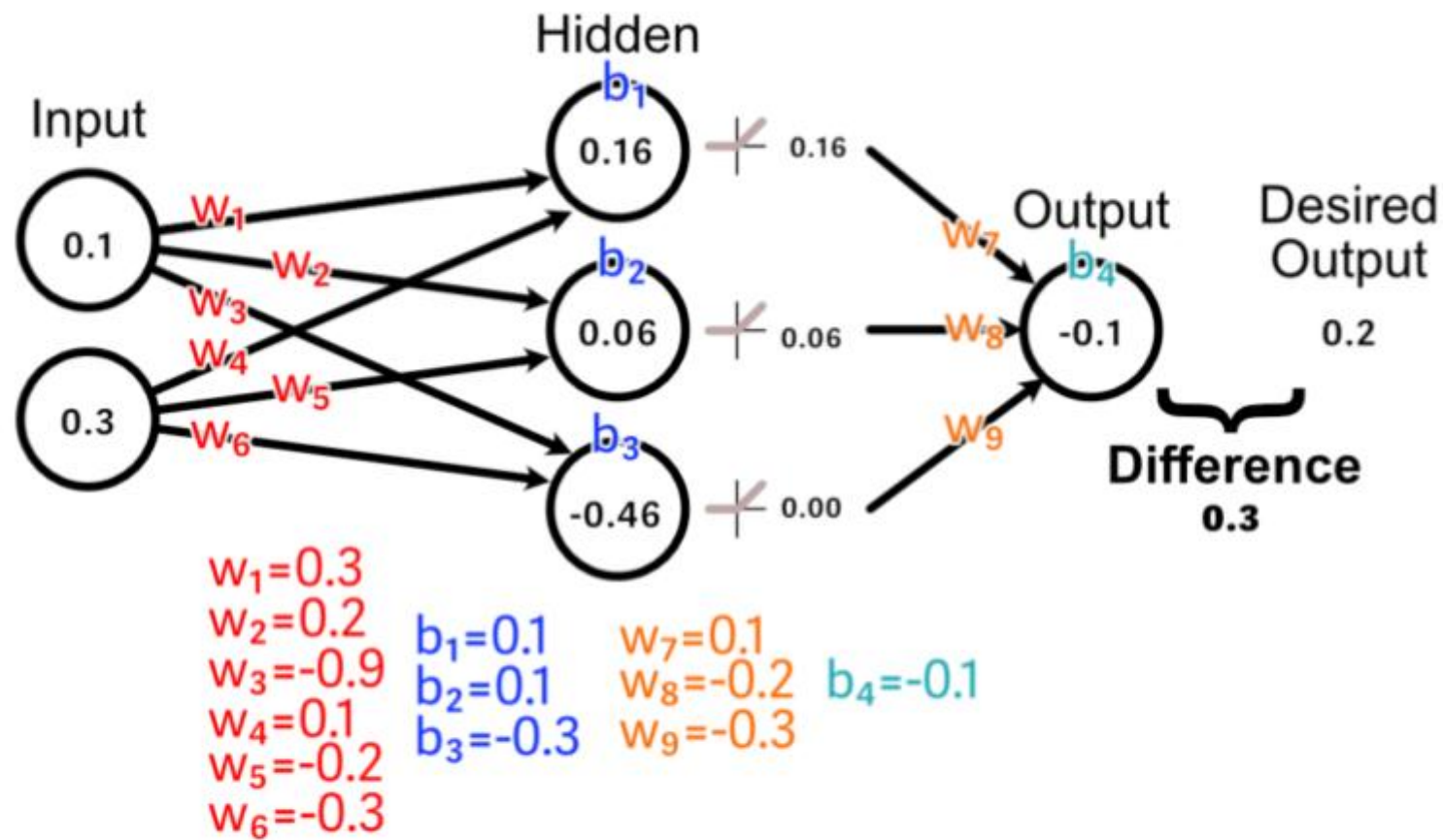
We'll apply the ReLU activation function to the value of our hidden perceptrons.

<b>Input</b>	<b>Desired Output</b>
0.1, 0.3	0.2
0.1, -0.1	0.0
-0.8, 0.2	-0.3
⋮	⋮

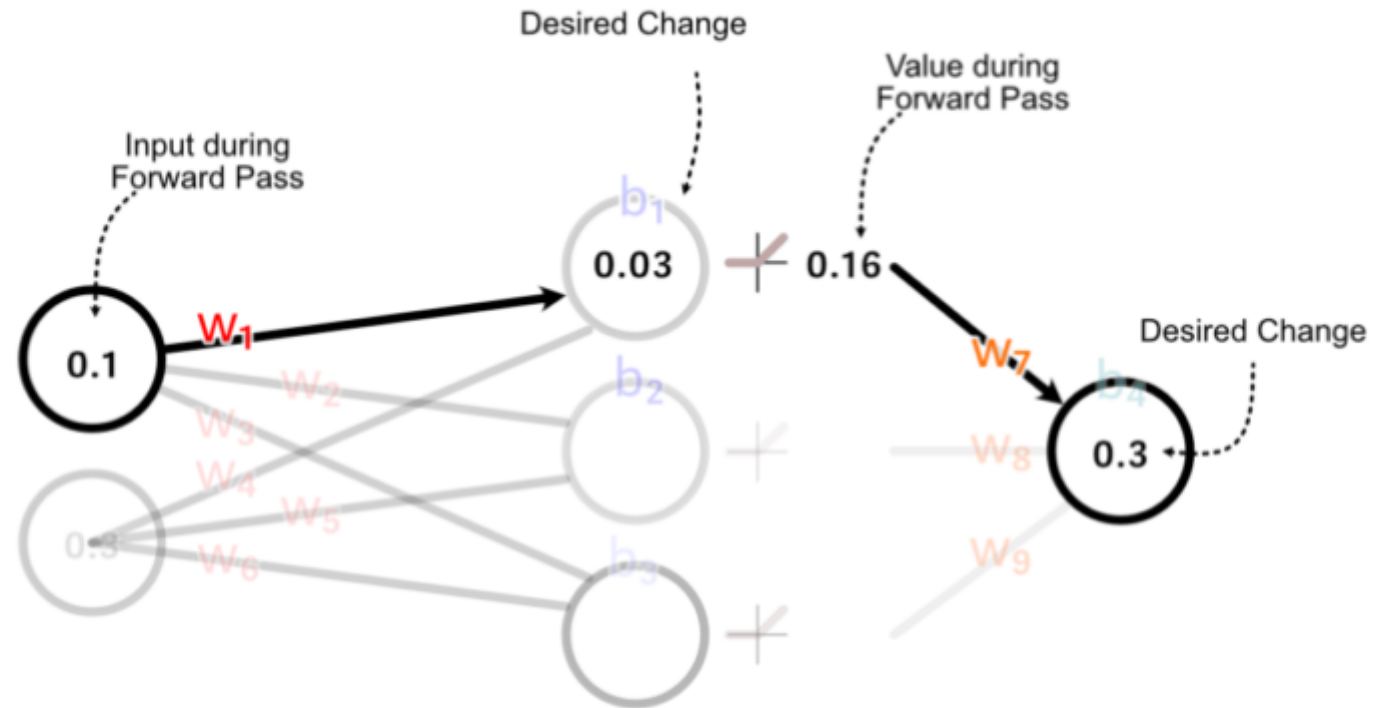
An example of the data that we'll be training off of.



Calculating the value of the hidden layer and output based on the input, including all major intermediary steps.

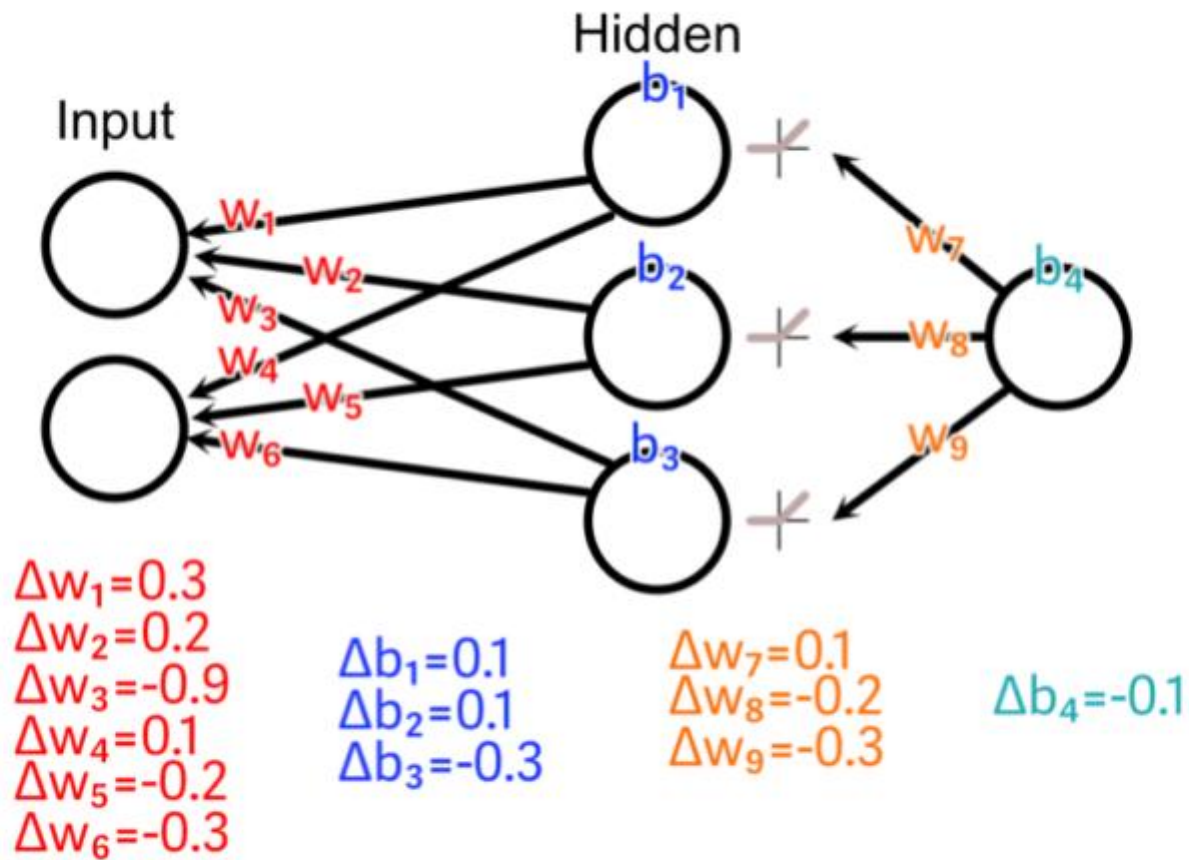


The training data has an input of 0.1 and 0.3, and the desired output (the average of the input) is 0.2. The prediction from the model was -0.1. Thus, the difference between the output and the desired output is 0.3.



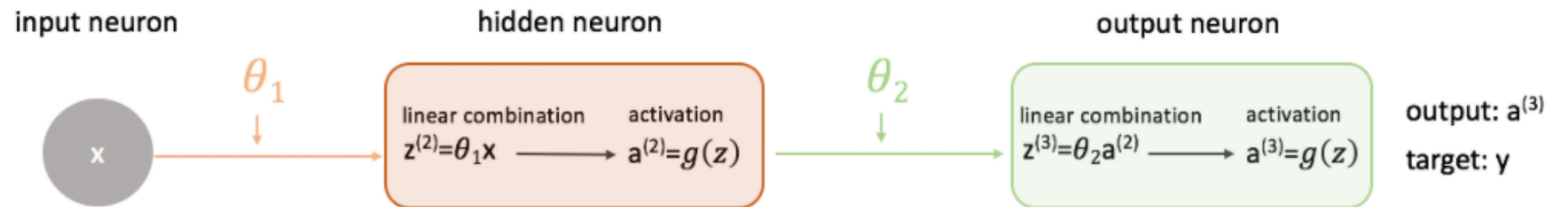
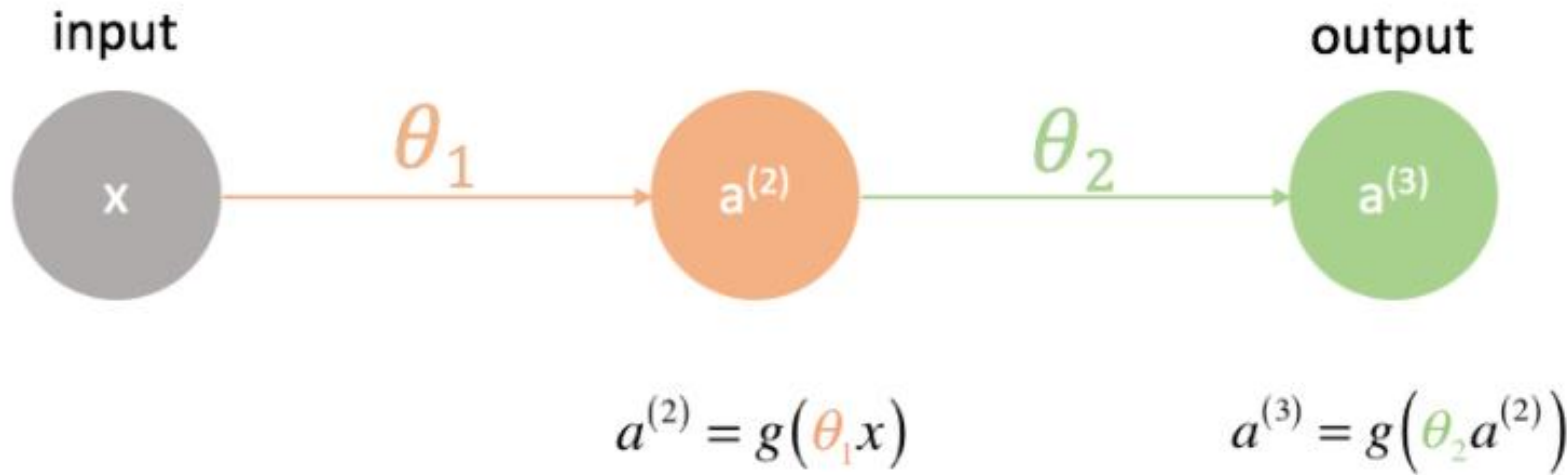
Change in Weight	=	Change in Output of Weight	x	Input To Weight
$\Delta w_7 = 0.050$	=	0.30	x	0.16
$\Delta w_1 = 0.003$	=	0.03	x	0.10

Now that we've calculated how the first hidden neuron should change, we can calculate how we should update  $w_1$  the same way we calculated how  $w_7$  should be updated previously.



By back propagating through the model, using a combination of values from the forward passes and desired changes from the backward pass at various points of the model, we can calculate how all parameters should

# Forward Propagation

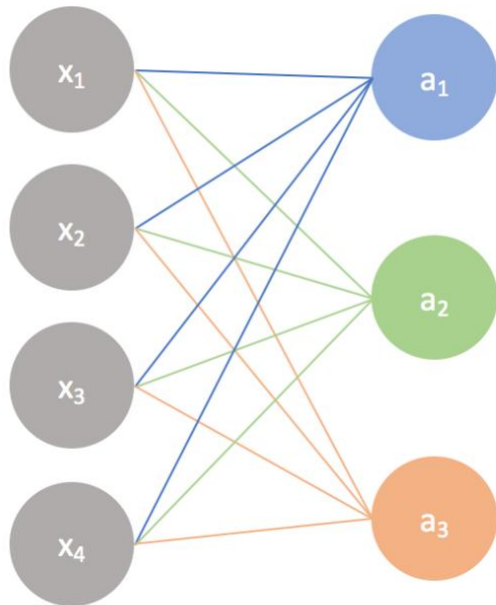




# Matrix Notation

Input layer

Output layer



$$\begin{bmatrix} w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 + b \\ w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 + b \\ w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 + b \end{bmatrix} \xrightarrow{\text{activation}} \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix}$$

$$A = \sigma(WX + B)$$



# *Cost Function ( $J$ )*

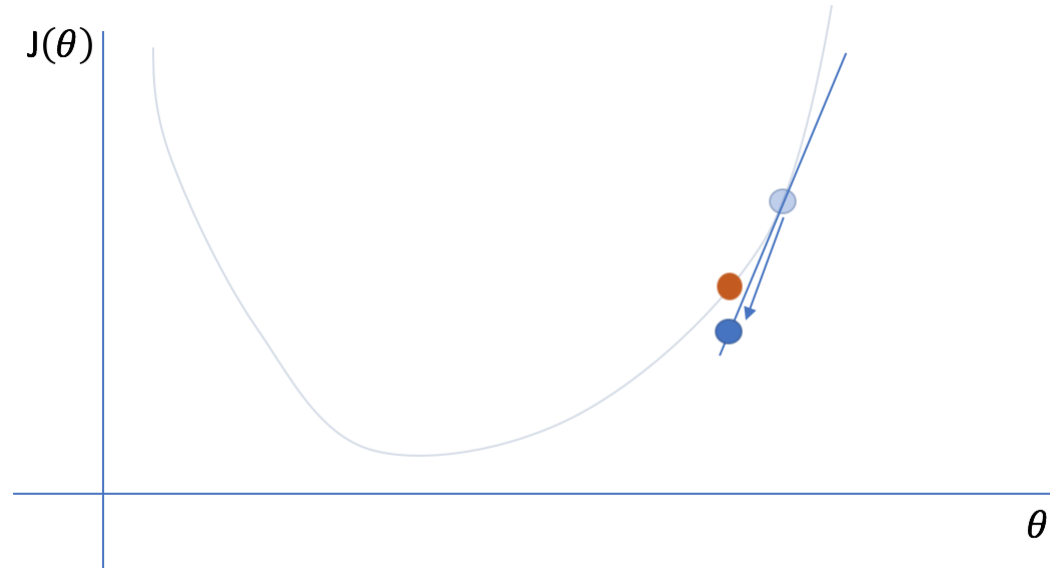
- Used to measure the error in a network prediction.
- Training a network: minimize the cost function.
- Examples:
  - Regression;
  - Mean Squared Error;
  - Mean Absolute Error;
  - Cross Entropy.

# Gradient Descent

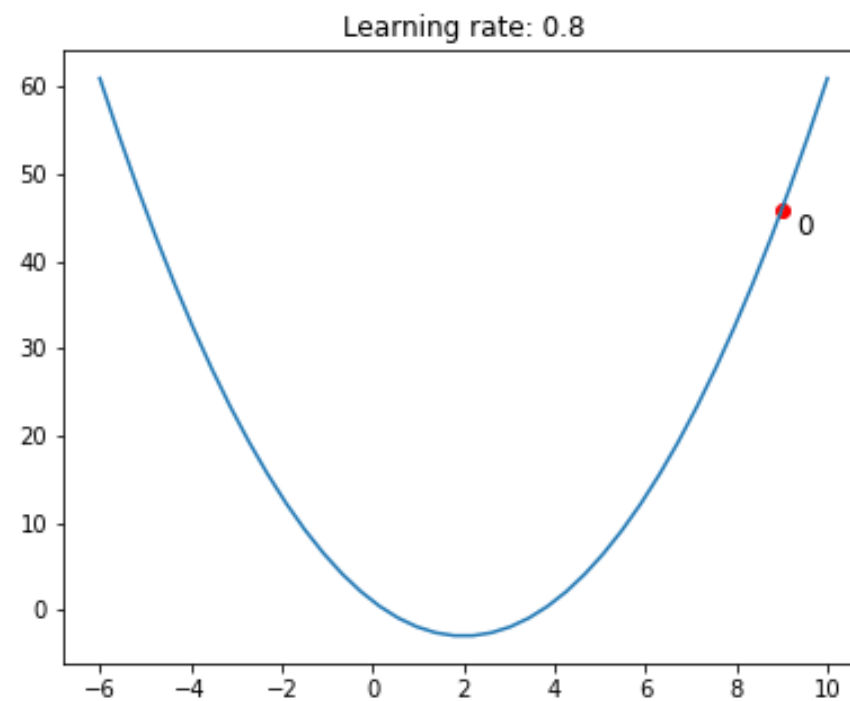
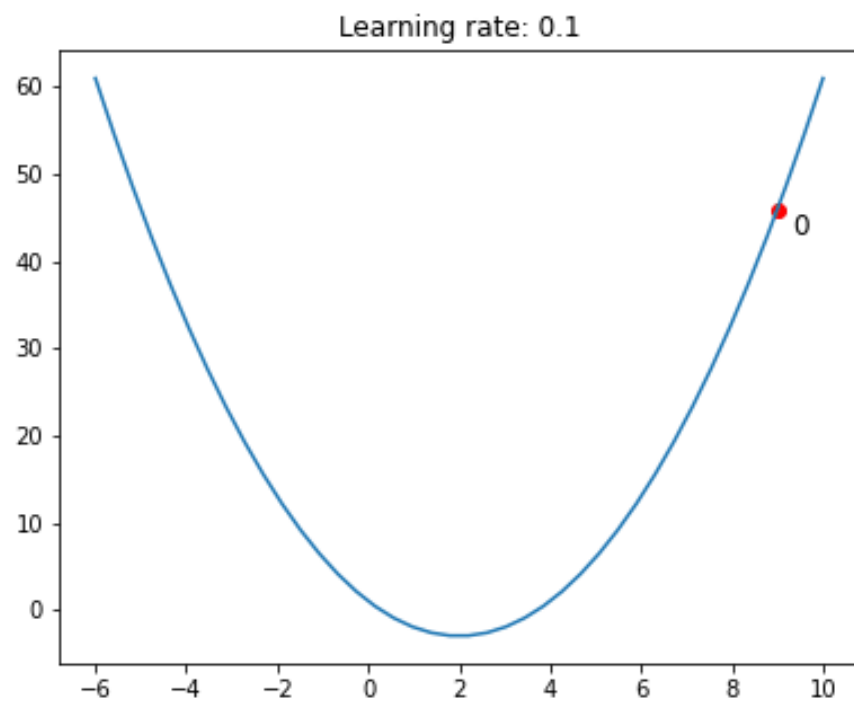
- Gradient: a vector that indicates the direction and magnitude in which, by displacement from the specified point, the greatest possible increase is obtained in the value of a quantity from which a scalar field for the space under consideration is defined.

$$\theta_i := \theta_i + \Delta\theta_i$$

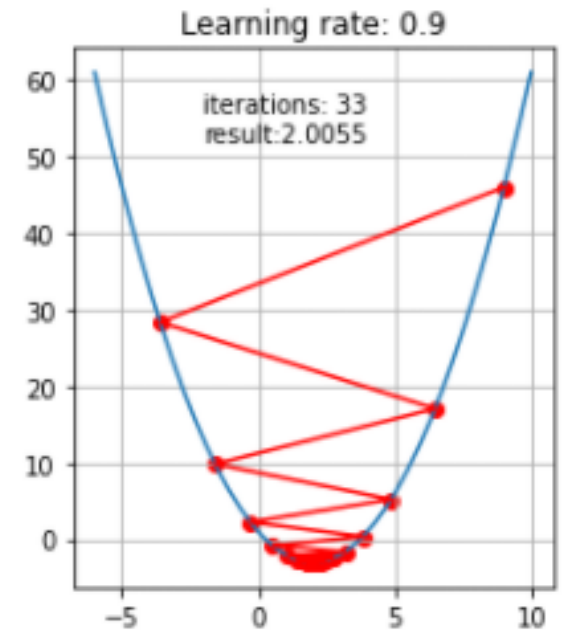
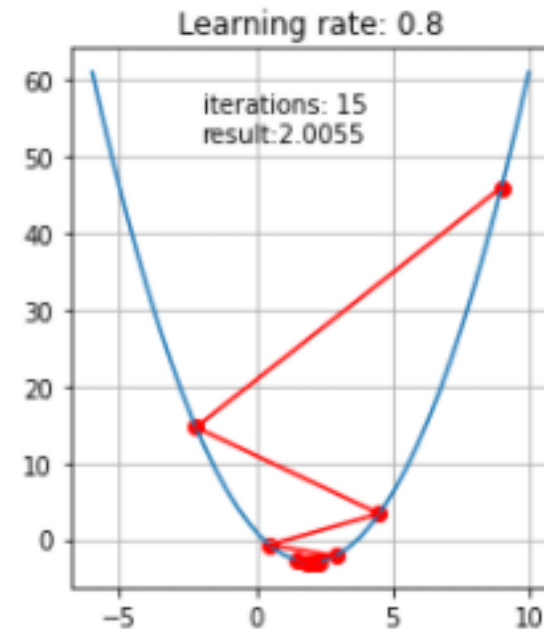
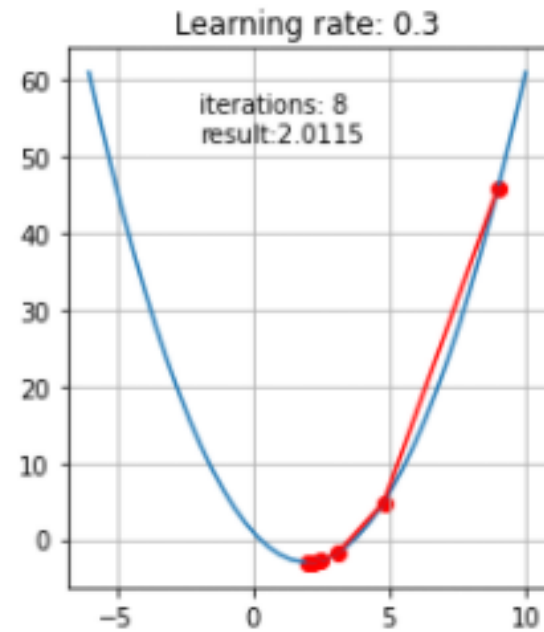
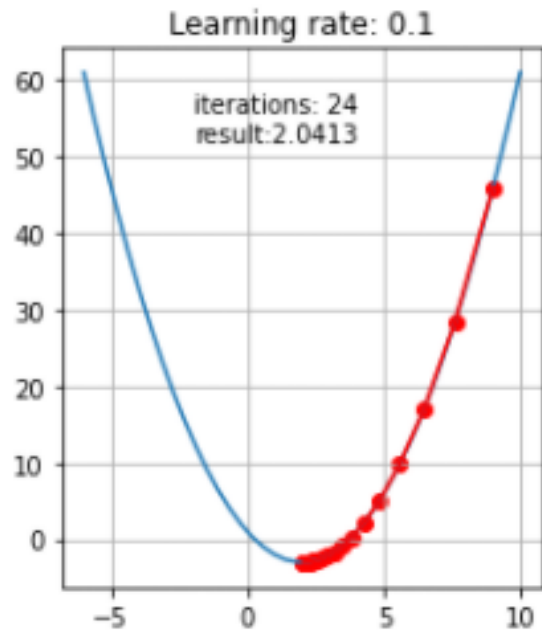
$$\Delta\theta_i = -\alpha \frac{\partial J(\theta)}{\partial \theta_i}$$



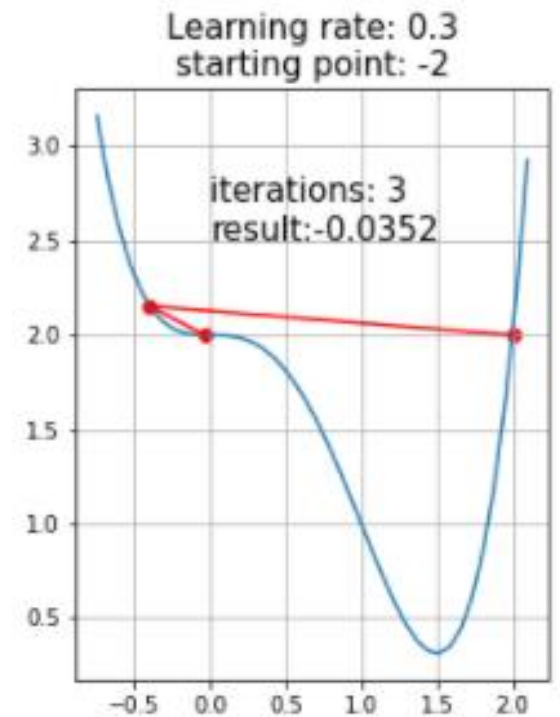
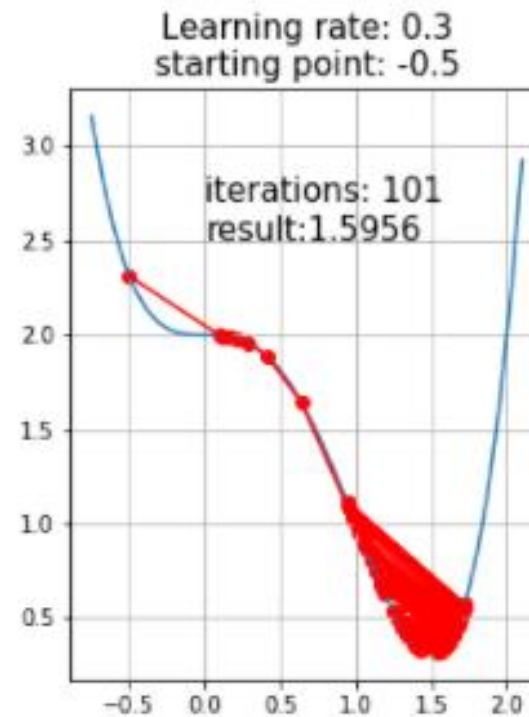
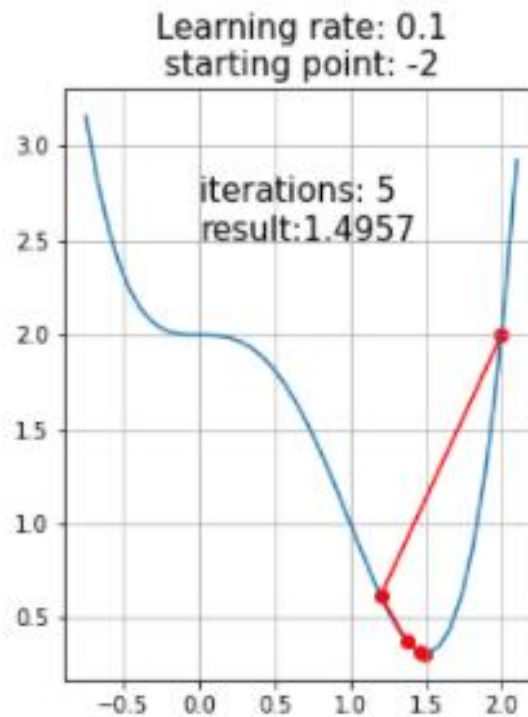
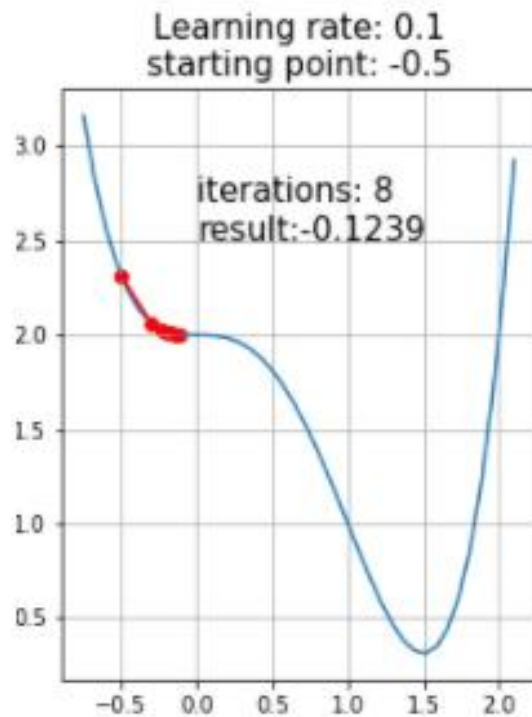
# *Learning Rate*



# *Different Learning rates*

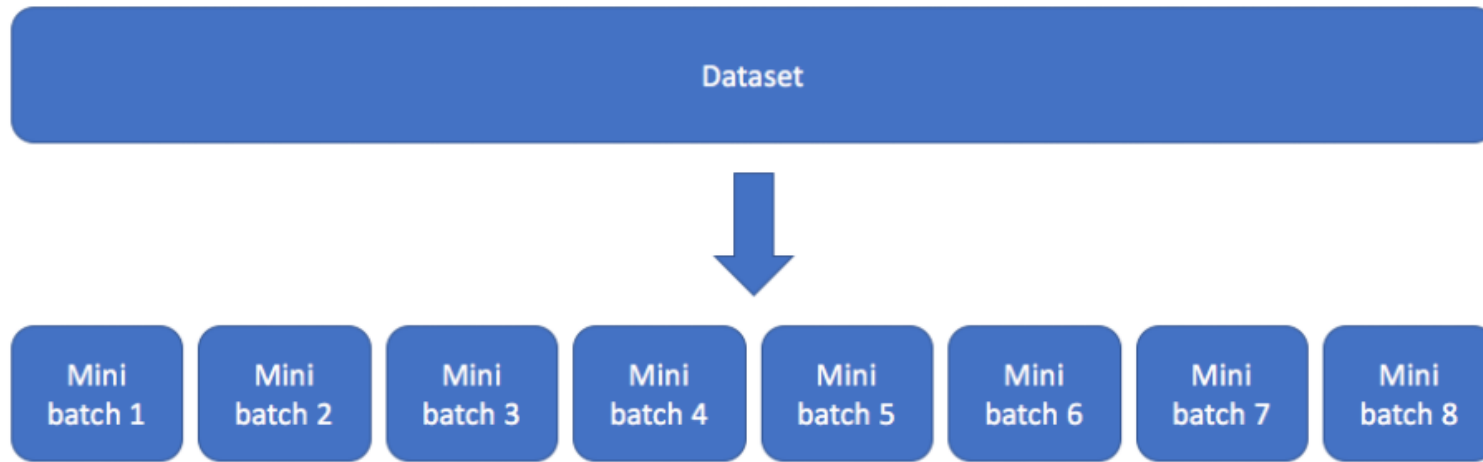


# *Different Learning rates*

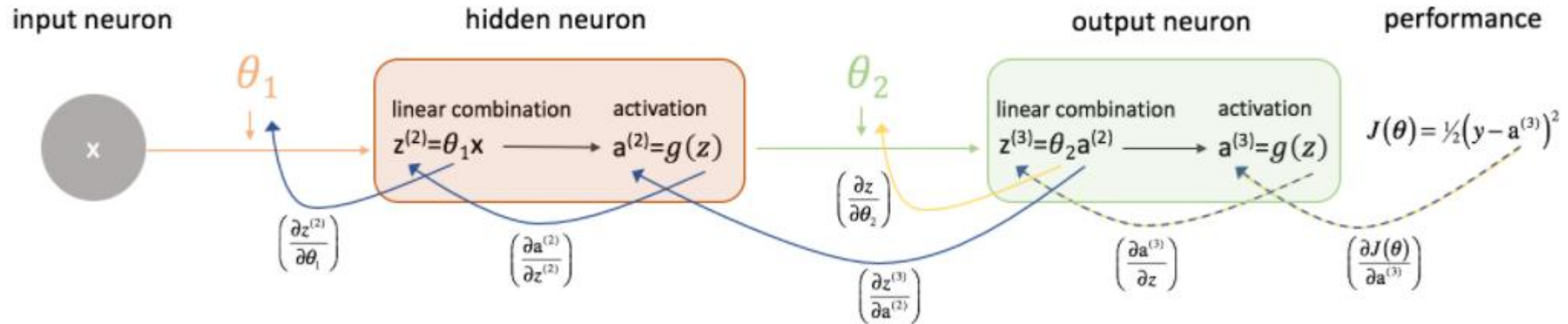


# Mini Batch

- Cost function analyzes the loss associated with each training example;
- The sum of the batch cost functions gives the overall cost of the entire dataset;
- Mini-batch: the process of "breaking down" the dataset into smaller chunks to calculate the cost function.



# Relationship between Weight and Cost Function



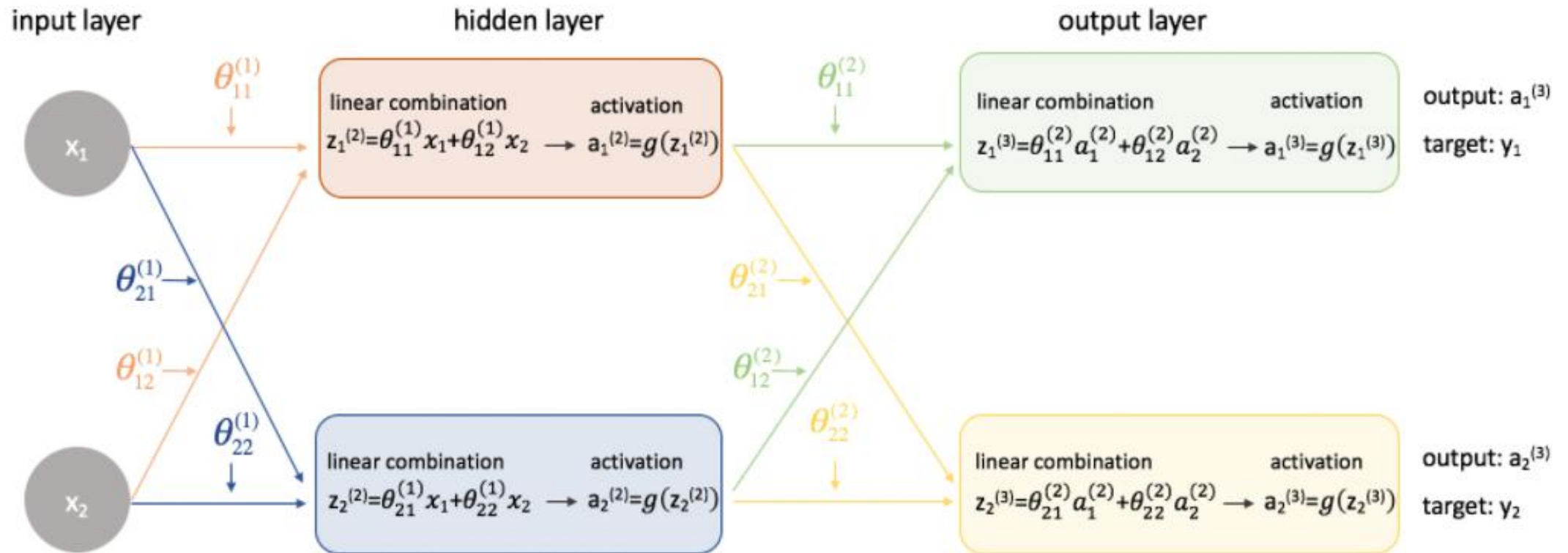
$$\frac{\partial J(\theta)}{\partial \theta_1} = ?$$

$$\frac{\partial J(\theta)}{\partial \theta_1} = \left( \frac{\partial J(\theta)}{\partial a^{(3)}} \right) \left( \frac{\partial a^{(3)}}{\partial z^{(3)}} \right) \left( \frac{\partial z^{(3)}}{\partial a^{(2)}} \right) \left( \frac{\partial a^{(2)}}{\partial z^{(2)}} \right) \left( \frac{\partial z^{(2)}}{\partial \theta_1} \right)$$

$$\frac{\partial J(\theta)}{\partial \theta_2} = ?$$

$$\frac{\partial J(\theta)}{\partial \theta_2} = \left( \frac{\partial J(\theta)}{\partial a^{(3)}} \right) \left( \frac{\partial a^{(3)}}{\partial z} \right) \left( \frac{\partial z}{\partial \theta_2} \right)$$

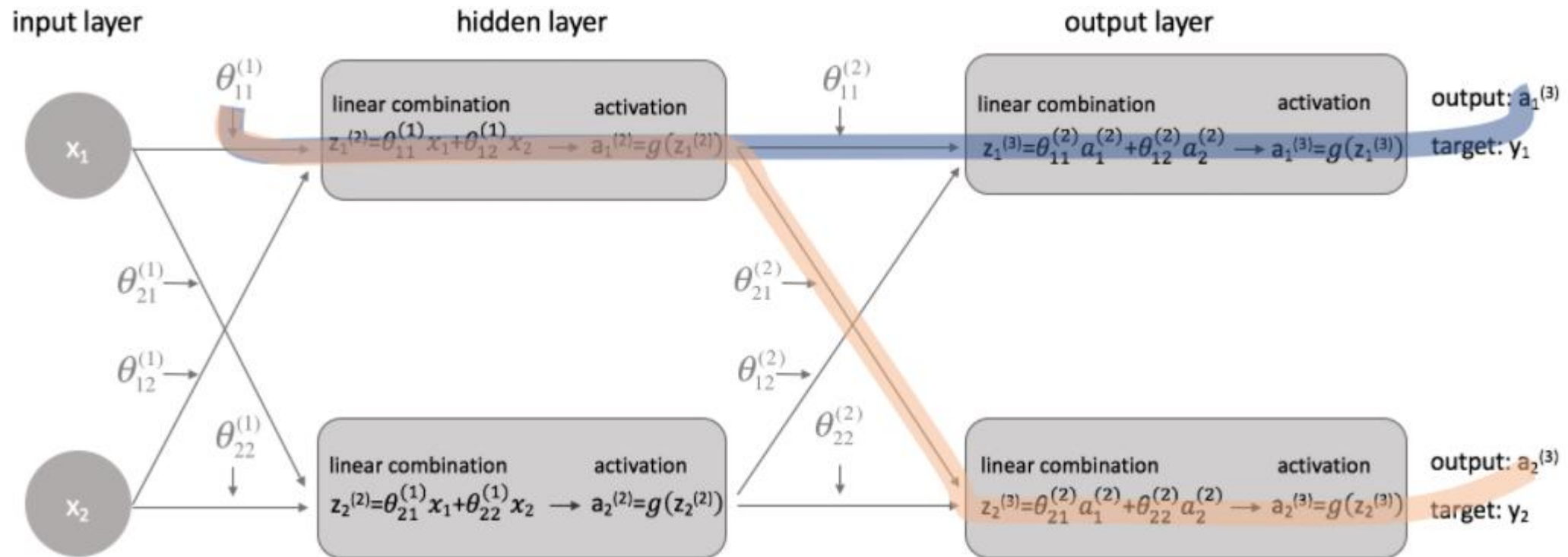
# Complicating the Problem: Two Layers and Two Outputs



$$J(\theta) = \frac{1}{2m} \sum (y_i - a_i^{(2)})^2$$



# Complicating the Problem: Two Layers and Two Outputs



$$\frac{\partial J(\theta)}{\partial \theta_{11}^{(1)}} = \left( \frac{\partial J(\theta)}{\partial a_1^{(3)}} \right) \left( \frac{\partial a_1^{(3)}}{\partial z_1^{(3)}} \right) \left( \frac{\partial z_1^{(3)}}{\partial a_1^{(2)}} \right) \left( \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} \right) \left( \frac{\partial z_1^{(2)}}{\partial \theta_{11}^{(1)}} \right) + \left( \frac{\partial J(\theta)}{\partial a_2^{(3)}} \right) \left( \frac{\partial a_2^{(3)}}{\partial z_2^{(3)}} \right) \left( \frac{\partial z_2^{(3)}}{\partial a_1^{(2)}} \right) \left( \frac{\partial a_1^{(2)}}{\partial z_1^{(2)}} \right) \left( \frac{\partial z_1^{(2)}}{\partial \theta_{11}^{(1)}} \right)$$



# *Complexity*

## 1 hidden layer and 1 output layer

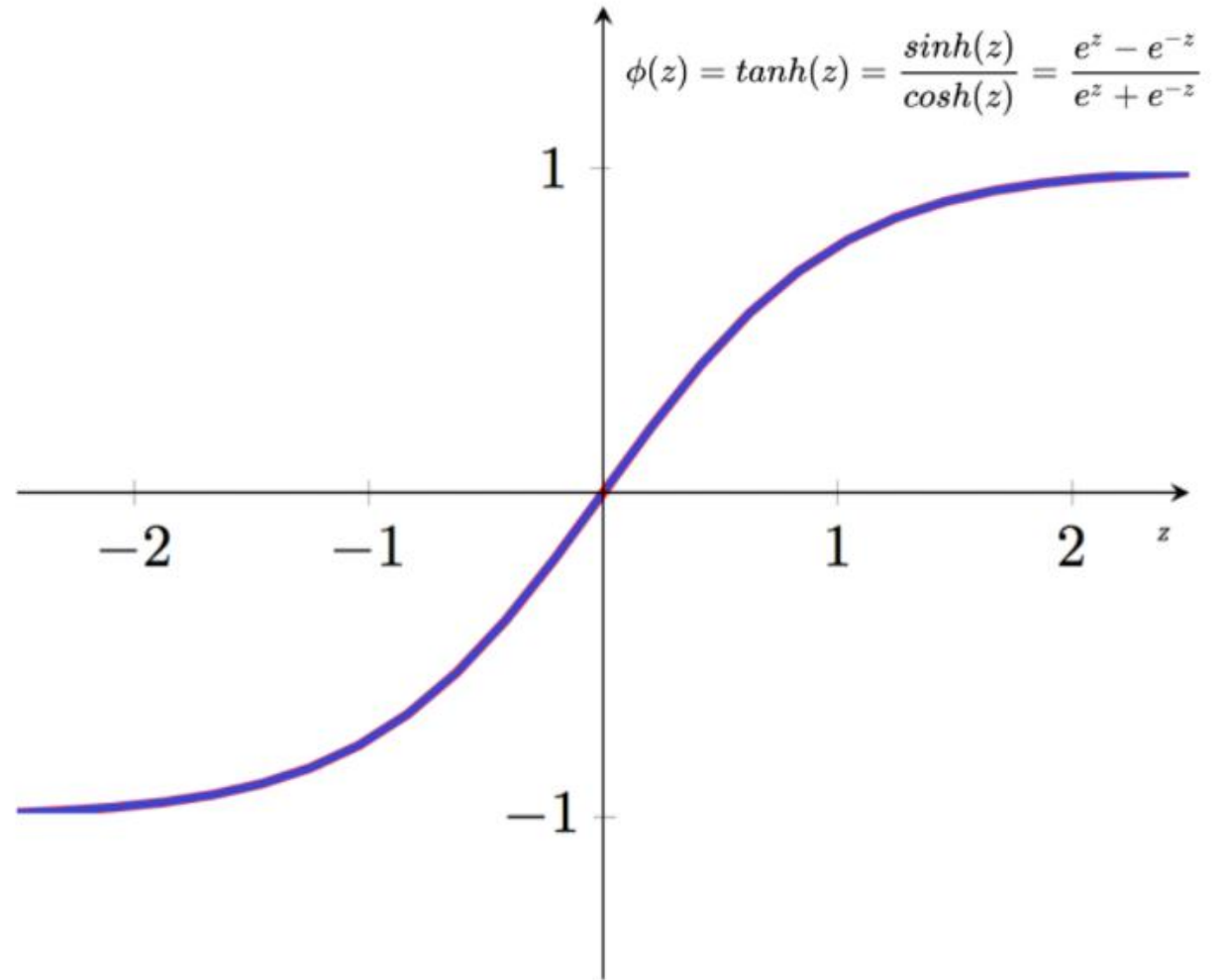
- 2 Parameters;
- 8 Partial Derivatives.

## 2 hidden layers and 2 output layers

- 8 Parameters;
- 52 Partial Derivatives.

# Activation Function

- Examples:
  - Boundary Functions;
  - Sigmoid Functions;
  - Rectified Linear Unit (ReLU) Functions;
  - Hyperbolic Tangent Functions.



*Thank You*



# Bias

- Allows shifting the activation function by adding a constant.

