

# Neural Networks

CISC 7026: Introduction to Deep Learning

University of Macau

1. We looked at linear and polynomial  $f$ 
  1. Looked at both classification and regression
  2. They have problems
    1. Input features scale poorly
    2. Bad performance around edges
  3. Neural networks fix many of these problems
  4. What is a neural network?
    1. Draw linear model as neural network
  5. Based on theory of the brain
    1. Invented ages ago
    2. Only recently have we learned to harness them

6. Neuron theory
  1. Connectivity
  2. Activation function
7. Parallels between real/artificial neuron
8. Matrix/graph duality
9. Single layer perceptron
10. Issues with one layer
  1. Not universal function approximator
11. Backprops
  1. Provides a way to train nn
    1. Assigns “fault” for each neuron

2. Recall closed form for linear model
  1. We use the gradient of the linear model
3. We use a similar approach

1. Limitations of linear models
2. History and overview of neural networks
3. Neurons
4. Perceptron
5. Multilayer Perceptron
6. Backpropagation
7. Gradient descent

We previously looked at linear and polynomial models for regression

We previously looked at linear and polynomial models for regression

$$f(\mathbf{x}, \boldsymbol{\theta}) = \theta_0 + \boldsymbol{\theta} \mathbf{x} = \theta_0 + \theta_1 x_1 + \theta_2 x_2, \dots$$

We previously looked at linear and polynomial models for regression

$$f(\mathbf{x}, \boldsymbol{\theta}) = \theta_0 + \boldsymbol{\theta} \mathbf{x} = \theta_0 + \theta_1 x_1 + \theta_2 x_2, \dots$$

$$\boldsymbol{\theta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$



Linear models are useful for simple problems

Linear models are useful for simple problems

Issues with very complex problems

Linear models are useful for simple problems

Issues with very complex problems

1. Poor scalability

Linear models are useful for simple problems

Issues with very complex problems

1. Poor scalability
2. Prone to overfitting

Linear models are useful for simple problems

Issues with very complex problems

1. Poor scalability
2. Prone to overfitting
3. Polynomials do not generalize well

Issues with very complex problems

1. **Poor scalability**
2. Polynomials do not generalize well

Polynomials fit tabular data well

Polynomials fit tabular data well

However, they scale poorly to higher-dimensional data like image pixels

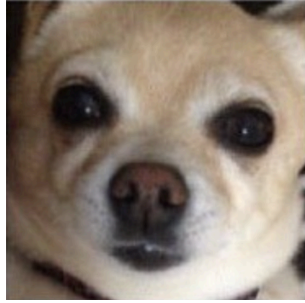


Polynomials fit tabular data well

However, they scale poorly to higher-dimensional data like image pixels

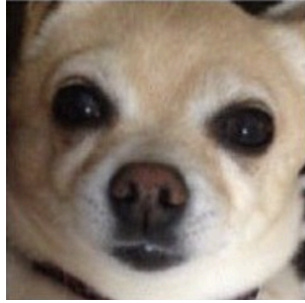


$$256 \times 256 \text{ pixels} = 65536 \text{ pixels}$$



$$256 \times 256 \text{ pixels} = 65536 \text{ pixels}$$

What does the design matrix look like for an n-degree polynomial?



$$256 \times 256 \text{ pixels} = 65536 \text{ pixels}$$

What does the design matrix look like for an n-degree polynomial?

$$\mathbf{X} = \begin{bmatrix} x_1^n & x_1^{n-1} & \dots & x_1^1 & 1 \\ x_2^n & x_2^{n-1} & \dots & x_2^1 & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x_p^n & x_p^{n-1} & \dots & x_p^1 & 1 \\ x_1^{n-1}x_2 & x_1^{n-2}x_2^2 & \dots & 0 & 1 \\ \vdots & \vdots & & \vdots & \vdots \end{bmatrix}$$

$$X = \begin{bmatrix} x_1^n & x_1^{n-1} & \dots & x_1^1 & 1 \\ x_2^n & x_2^{n-1} & \dots & x_2^1 & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x_p^n & x_p^{n-1} & \dots & x_p^1 & 1 \\ x_1^{n-1}x_2 & x_1^{n-2}x_2^2 & \dots & 0 & 1 \\ \vdots & \vdots & & \vdots & \vdots \end{bmatrix}$$

**Question:** How big is the matrix for 65,536 pixels and  $n = 3$ ?

**Answer:**  $65,536^3 \approx 10^{14}$  parameters

$$X = \begin{bmatrix} x_1^n & x_1^{n-1} & \dots & x_1^1 & 1 \\ x_2^n & x_2^{n-1} & \dots & x_2^1 & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x_p^n & x_p^{n-1} & \dots & x_p^1 & 1 \\ x_1^{n-1}x_2 & x_1^{n-2}x_2^2 & \dots & 0 & 1 \\ \vdots & \vdots & & \vdots & \vdots \end{bmatrix}$$

**Question:** How big is the matrix for 65,536 pixels and  $n = 3$ ?

**Answer:**  $65,536^3 \approx 10^{14}$  parameters

**Question:** How big is the matrix for 65,536 pixels and  $n = 3$ ?

**Answer:**  $65,536^3 \approx 10^{14}$  parameters

**Question:** How big is the matrix for 65,536 pixels and  $n = 3$ ?

**Answer:**  $65,536^3 \approx 10^{14}$  parameters

We must invert  $\mathbf{X}^\top \mathbf{X}$ , requiring  $O(n^3)$  time



**Question:** How big is the matrix for 65,536 pixels and  $n = 3$ ?

**Answer:**  $65,536^3 \approx 10^{14}$  parameters

We must invert  $\mathbf{X}^\top \mathbf{X}$ , requiring  $O(n^3)$  time

Largest matrix ever inverted is  $\approx 10^{12}$

**Question:** How big is the matrix for 65,536 pixels and  $n = 3$ ?

**Answer:**  $65,536^3 \approx 10^{14}$  parameters

We must invert  $\mathbf{X}^\top \mathbf{X}$ , requiring  $O(n^3)$  time

Largest matrix ever inverted is  $\approx 10^{12}$

For comparison, GPT-4 has  $10^{12}$  parameters

**Question:** How big is the matrix for 65,536 pixels and  $n = 3$ ?

**Answer:**  $65,536^3 \approx 10^{14}$  parameters

We must invert  $\mathbf{X}^\top \mathbf{X}$ , requiring  $O(n^3)$  time

Largest matrix ever inverted is  $\approx 10^{12}$

For comparison, GPT-4 has  $10^{12}$  parameters

Polynomial regression scales poorly to high dimensional data

Issues with very complex problems

1. **Poor scalability**
2. Polynomials do not generalize well

Issues with very complex problems

1. Poor scalability
2. **Polynomials do not generalize well**

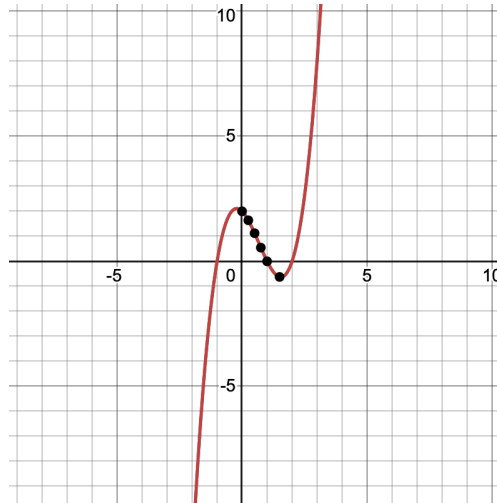
Polynomials tend towards  $-\infty, \infty$  outside of the support

Polynomials tend towards  $-\infty, \infty$  outside of the support

$$f(x) = x^3 - 2x^2 - x + 2$$

Polynomials tend towards  $-\infty, \infty$  outside of the support

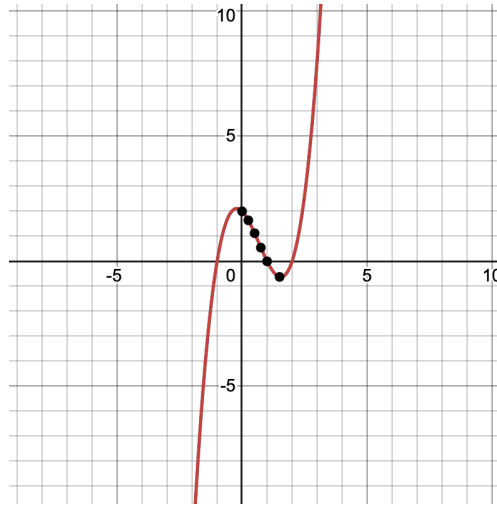
$$f(x) = x^3 - 2x^2 - x + 2$$





Polynomials tend towards  $-\infty, \infty$  outside of the support

$$f(x) = x^3 - 2x^2 - x + 2$$



If breed of dog missing from training set, we still want to classify it as dog!

Linear and polynomial regression have issues

# Linear and polynomial regression have issues

1. Poor scalability

Linear and polynomial regression have issues

1. Poor scalability
2. Polynomials do not generalize well

Relax

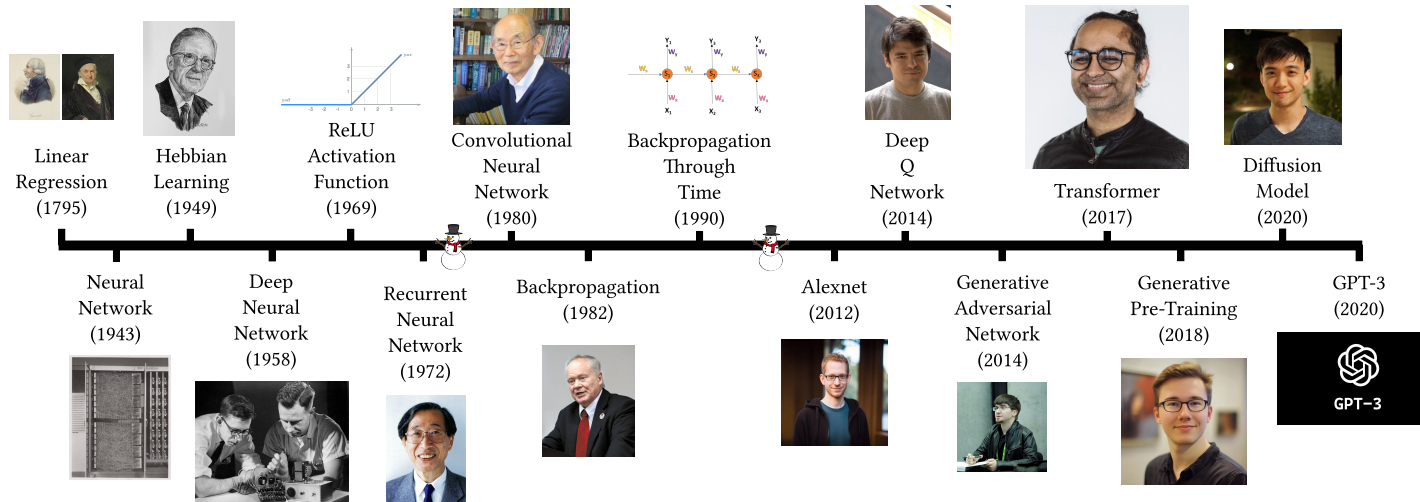
Can we improve upon the linear/polynomial model?

Can we improve upon the linear/polynomial model?

Yes, with neural networks

# Can we improve upon the linear/polynomial model?

## Yes, with neural networks





**Brain:** Biological neurons  $\rightarrow$  Biological neural network

**Brain:** Biological neurons  $\rightarrow$  Biological neural network

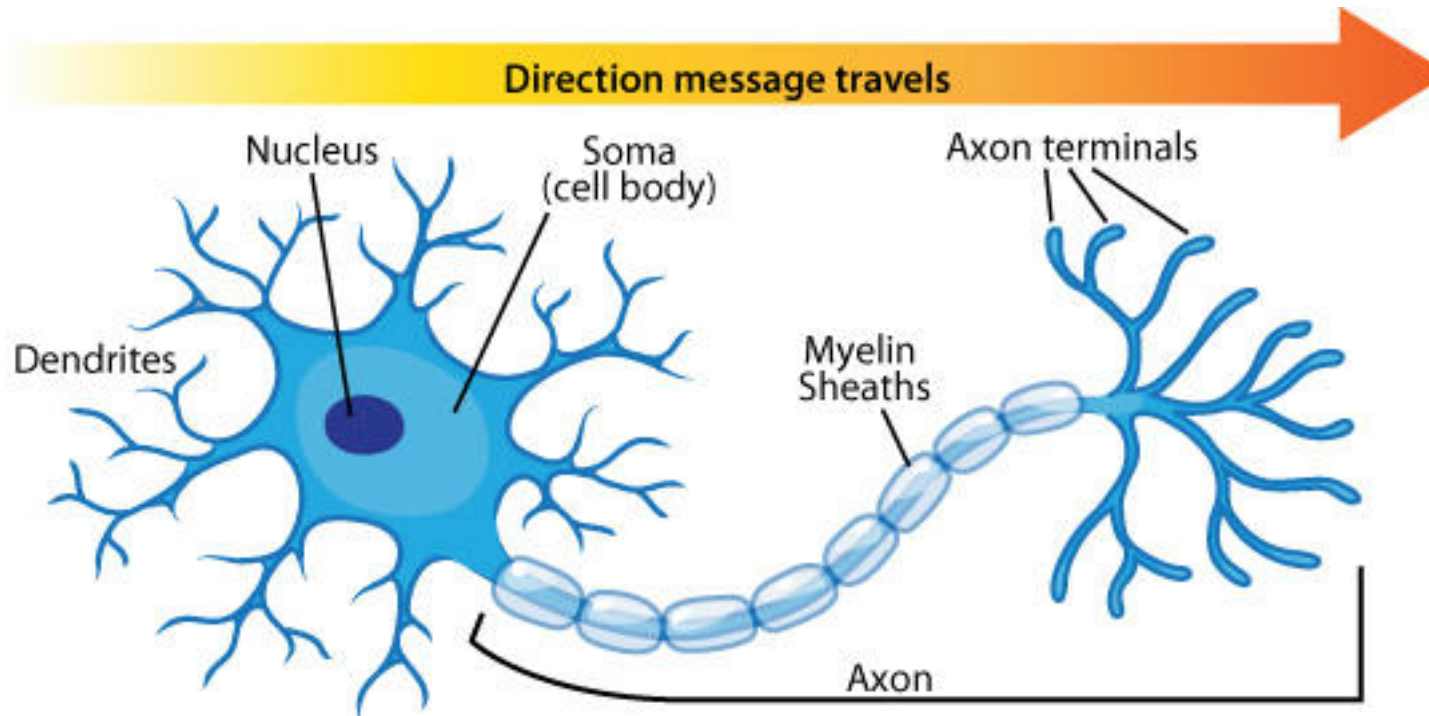
**Computer:** Artificial neurons  $\rightarrow$  Artificial neural network

**Brain:** Biological neurons  $\rightarrow$  Biological neural network

**Computer:** Artificial neurons  $\rightarrow$  Artificial neural network

Neurons send and receive electrical impulses along axons and dendrites

Neurons send and receive electrical impulses along axons and dendrites



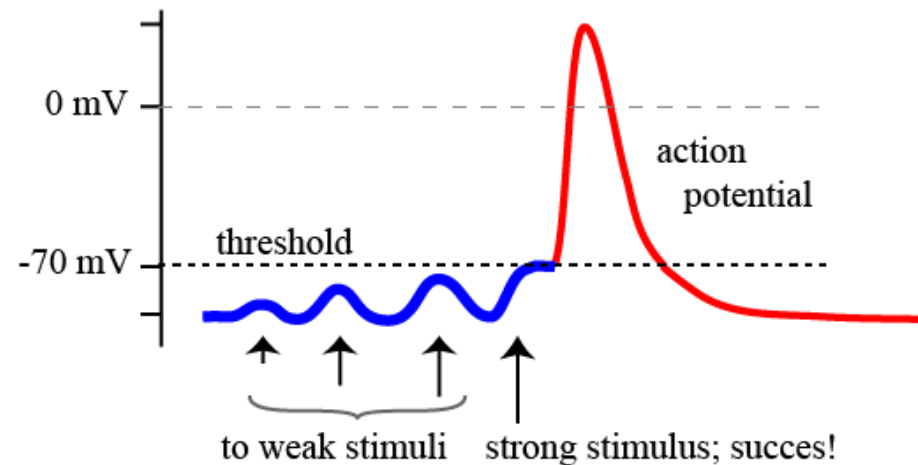
How does a neuron send an impulse (“fire”)?

How does a neuron send an impulse (“fire”)?

Incoming impulses (via dendrites) change the electric potential of the neuron

How does a neuron send an impulse (“fire”)?

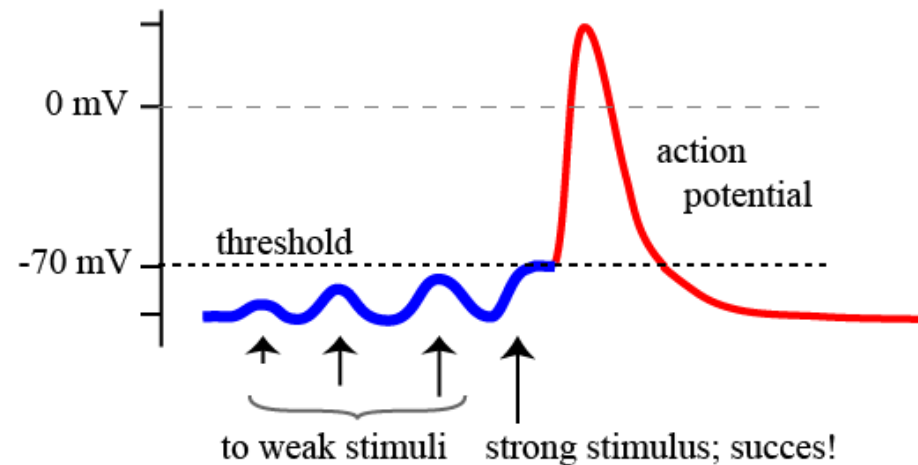
Incoming impulses (via dendrites) change the electric potential of the neuron



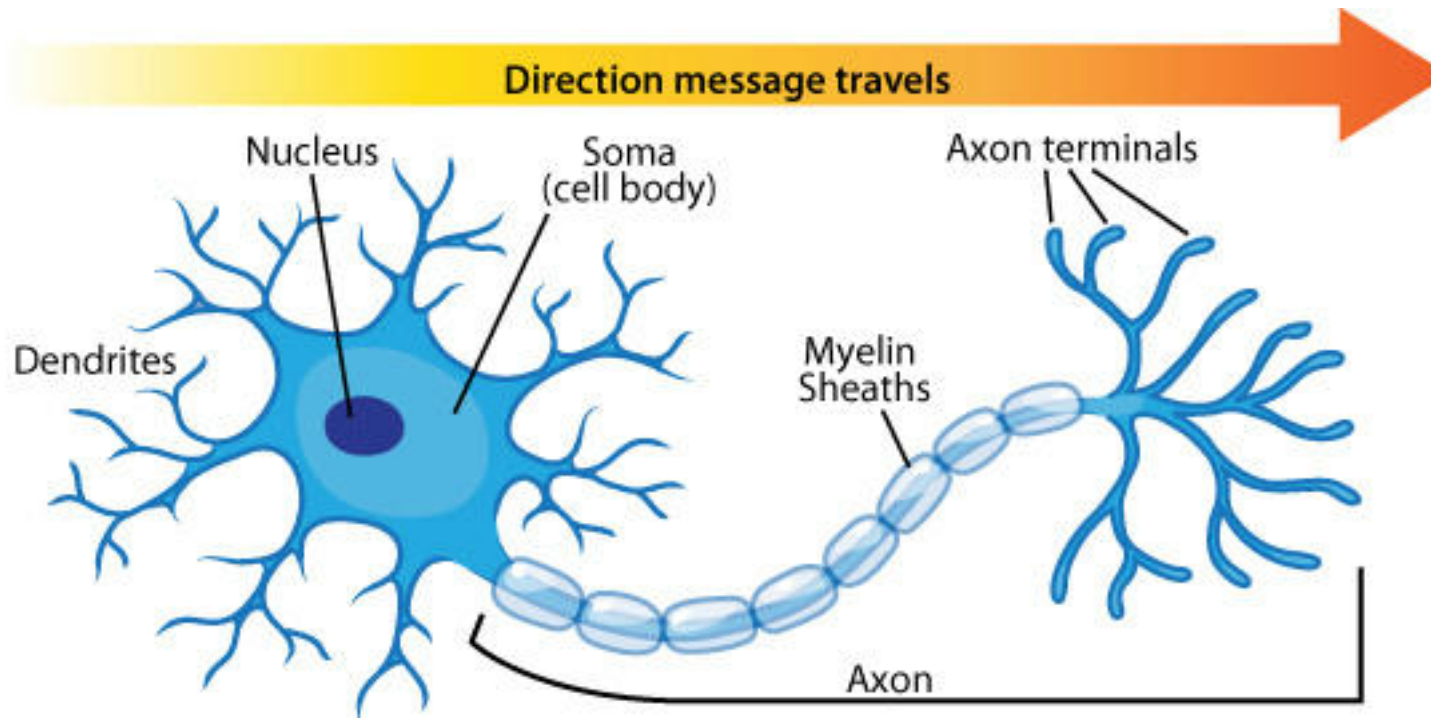


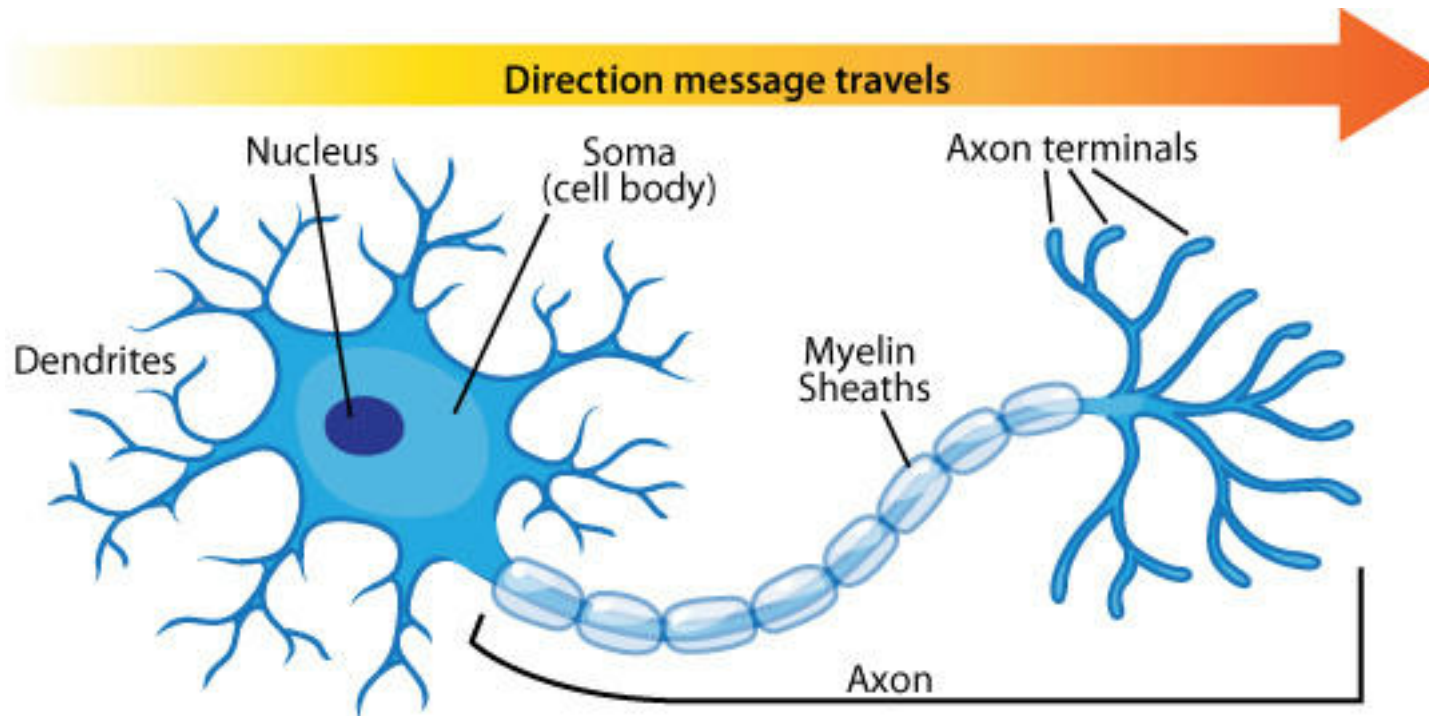
How does a neuron send an impulse (“fire”)?

Incoming impulses (via dendrites) change the electric potential of the neuron



Pain triggers initial nerve impulse, sets of impulse chain into the brain

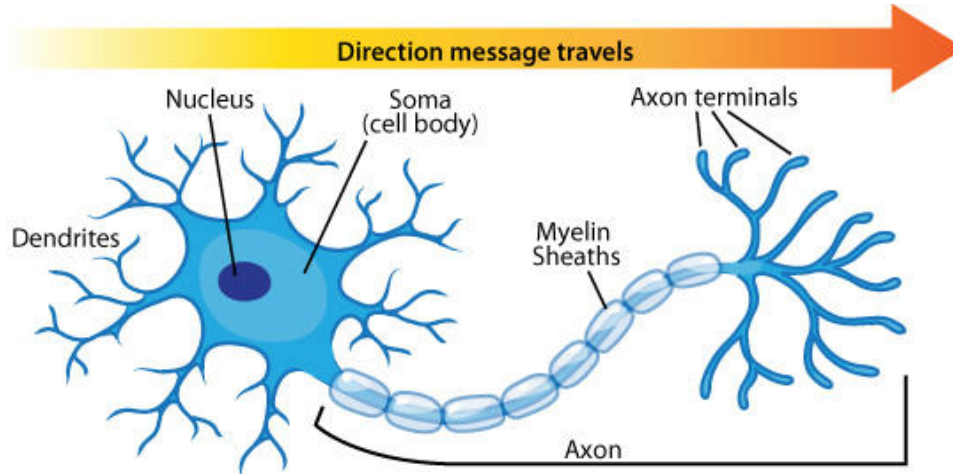




**Question:** How would you model a neuron mathematically?

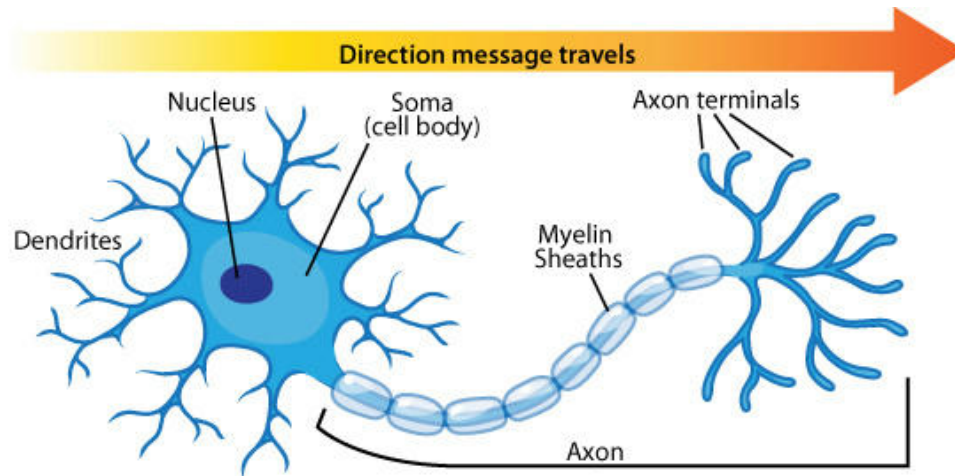
Let us define a neuron as a function

Let us define a neuron as a function



Neuron has a structure of  
dendrites

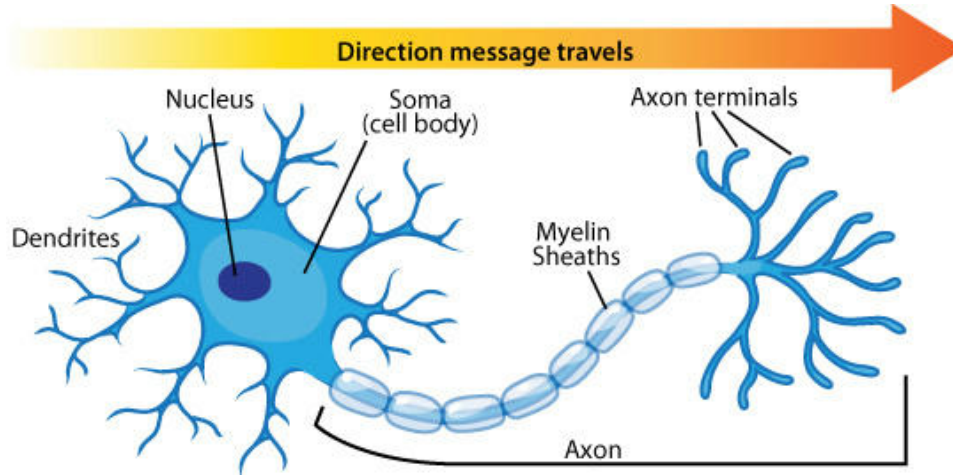
Let us define a neuron as a function



Neuron has a structure of  
dendrites

$$f \left( \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \right) = f \left( \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 1 \end{bmatrix} \right)$$

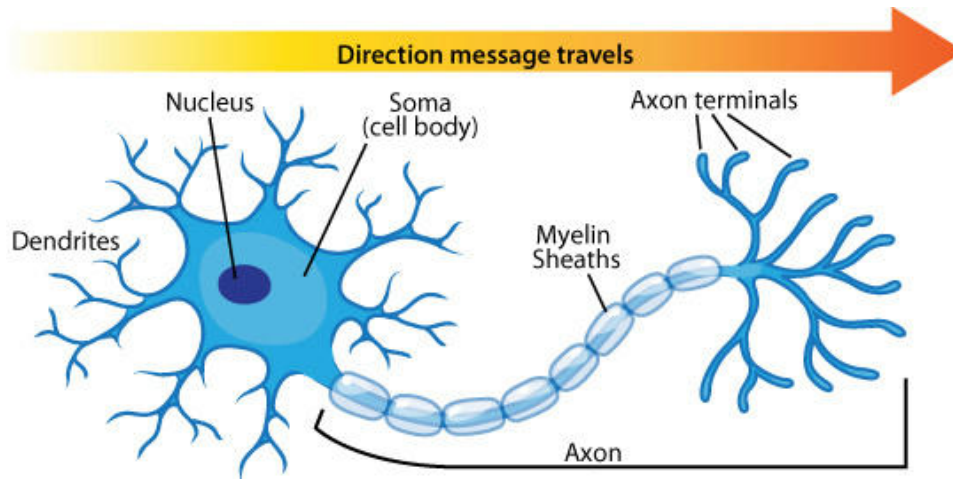
Let us define a neuron as a function



Each incoming dendrite has some voltage potential

Let us define a neuron as a function

Each incoming dendrite has some voltage potential

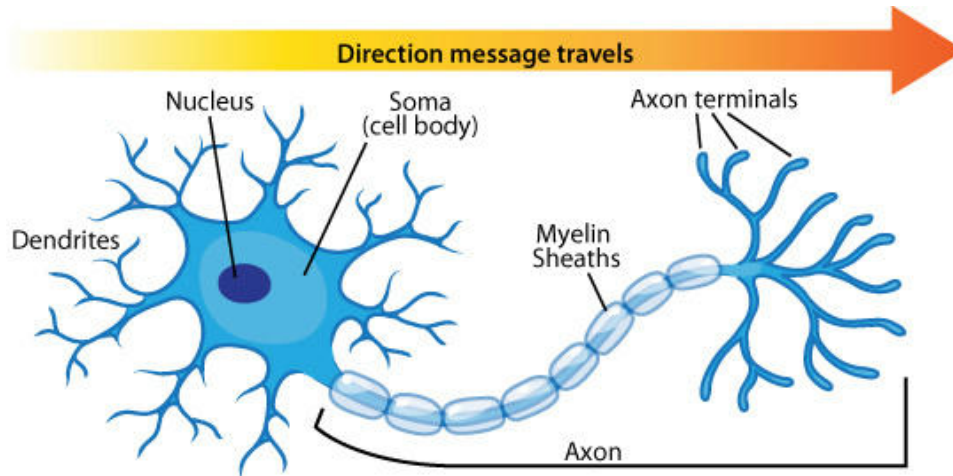


$$f\left(\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}\right)$$

$$\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} 0.5 \\ \vdots \\ -0.3 \end{bmatrix}$$

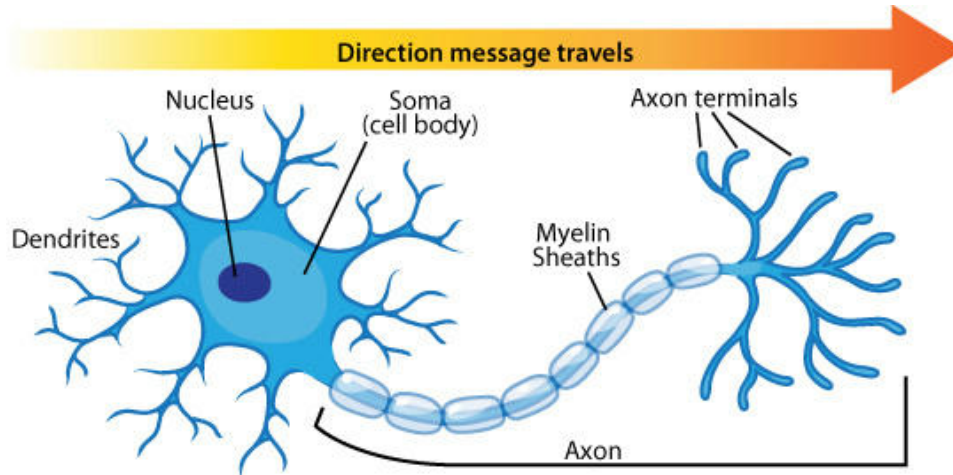


Let us define a neuron as a function



Voltage potentials sum together to give us the voltage in the cell body

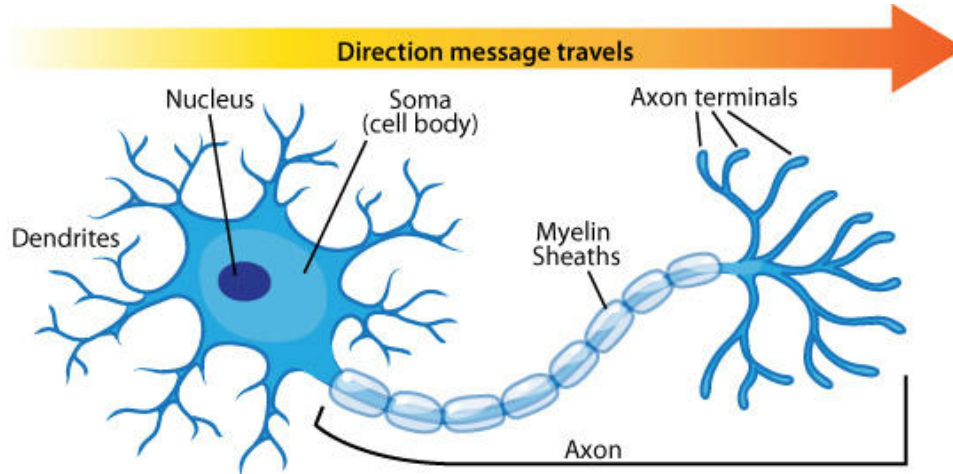
Let us define a neuron as a function



Voltage potentials sum together to give us the voltage in the cell body

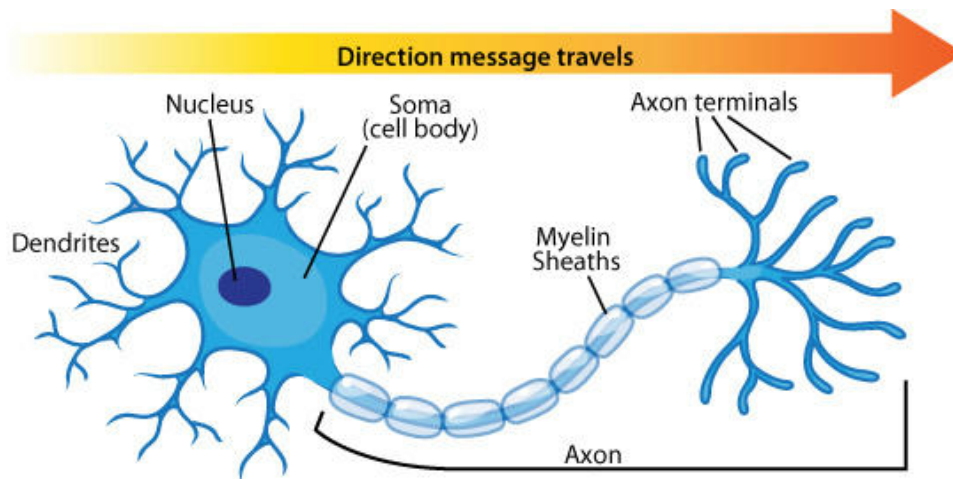
$$f\left(\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}\right) = \sum_{i=1}^n x_i \theta_i$$

Let us define a neuron as a function



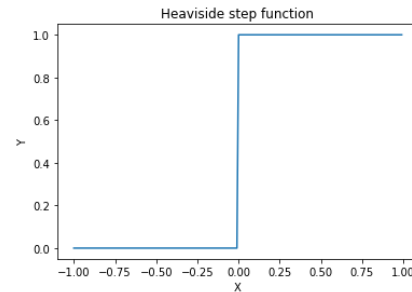
The axon fires only if the voltage is over a threshold

Let us define a neuron as a function



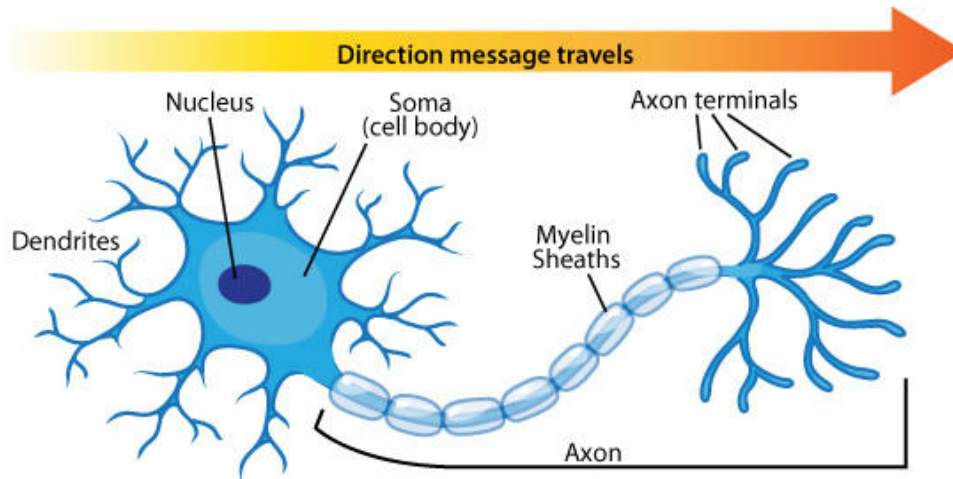
The axon fires only if the voltage is over a threshold

$$\sigma(x) =$$

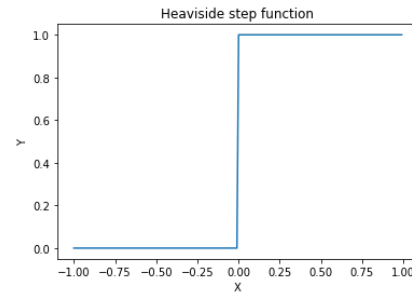


Let us define a neuron as a function

The axon fires only if the voltage is over a threshold



$$\sigma(x) =$$



$$f\left(\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}\right) = \sigma\left(\sum_{i=1}^n x_i \theta_i\right)$$

This is almost the artificial neuron!

$$f\left(\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}\right) = \sigma\left(\sum_{i=1}^n x_i \theta_i\right)$$

This is almost the artificial neuron!

$$f\left(\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}\right) = \sigma\left(\sum_{i=1}^n x_i \theta_i\right)$$

$$f(\mathbf{x}, \boldsymbol{\theta}) = \sigma\left(\sum_{i=1}^n x_i \theta_i\right)$$

**Question:** Does it look familiar to any other functions we have seen?

This is almost the artificial neuron!

$$f\left(\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}\right) = \sigma\left(\sum_{i=1}^n x_i \theta_i\right)$$

$$f(\mathbf{x}, \boldsymbol{\theta}) = \sigma\left(\sum_{i=1}^n x_i \theta_i\right)$$

**Question:** Does it look familiar to any other functions we have seen?

**Answer:** The linear model!



$$f(\boldsymbol{x}, \boldsymbol{\theta}) = \sigma \left( \sum_{i=1}^n x_i \theta_i \right)$$

Artificial neuron

$$f(\mathbf{x}, \boldsymbol{\theta}) = \sigma \left( \sum_{i=1}^n x_i \theta_i \right)$$

Artificial neuron

$$f(\mathbf{x}, \boldsymbol{\theta}) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

Linear model

$$f(\mathbf{x}, \boldsymbol{\theta}) = \sigma \left( \sum_{i=1}^n x_i \theta_i \right)$$

Artificial neuron

$$f(\mathbf{x}, \boldsymbol{\theta}) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

Linear model

It is the linear model with an activation function!

$$f(\mathbf{x}, \boldsymbol{\theta}) = \sigma \left( \sum_{i=1}^n x_i \theta_i \right)$$

Artificial neuron

$$f(\mathbf{x}, \boldsymbol{\theta}) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

Linear model

It is the linear model with an activation function!

We add a bias term to the neuron, for the same reason we add a bias term to the linear model

$$f(\mathbf{x}, \boldsymbol{\theta}) = \sigma \left( \sum_{i=1}^n x_i \theta_i \right) \quad \text{Artificial neuron}$$

$$f(\mathbf{x}, \boldsymbol{\theta}) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n \quad \text{Linear model}$$

It is the linear model with an activation function!

We add a bias term to the neuron, for the same reason we add a bias term to the linear model

$$f(\mathbf{x}, \boldsymbol{\theta}) = \sigma \left( \theta_0 + \sum_{i=1}^n x_i \theta_i \right)$$

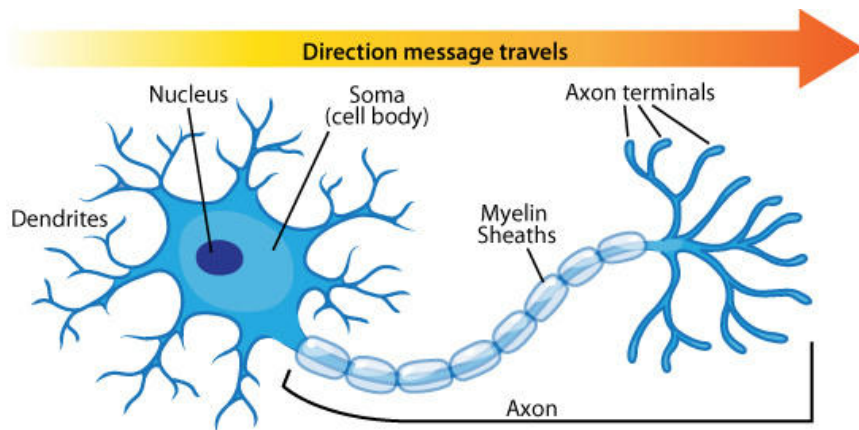
$$f(\mathbf{x}, \boldsymbol{\theta}) = \sigma \left( \sum_{i=1}^n x_i \theta_i \right) \quad \text{Artificial neuron}$$

$$f(\mathbf{x}, \boldsymbol{\theta}) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n \quad \text{Linear model}$$

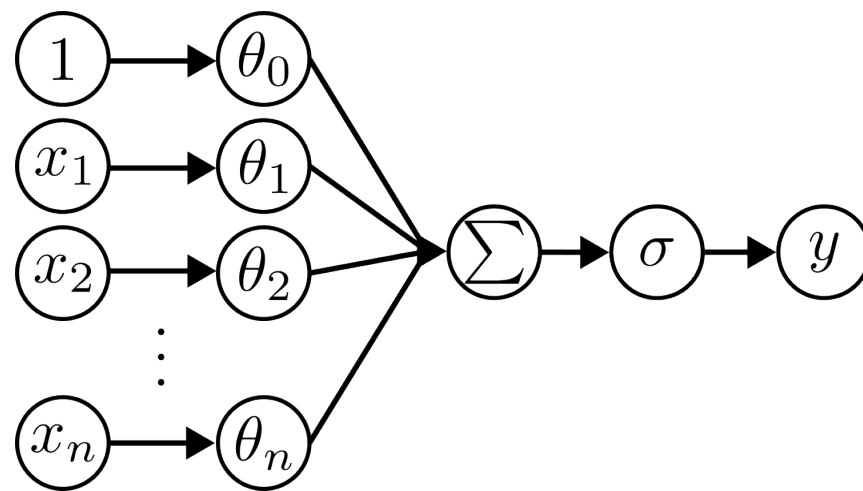
It is the linear model with an activation function!

We add a bias term to the neuron, for the same reason we add a bias term to the linear model

$$f(\mathbf{x}, \boldsymbol{\theta}) = \sigma \left( \theta_0 + \sum_{i=1}^n x_i \theta_i \right)$$



$$f\left(\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_n \end{bmatrix}\right) = \sigma\left(\theta_0 + \sum_{i=1}^n x_i \theta_i\right)$$



We can also write a neuron in terms of a dot product



We can also write a neuron in terms of a dot product

$$f\left(\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_n \end{bmatrix}\right) = \sigma(\theta_0 + \boldsymbol{\theta}_{1:n} \cdot \boldsymbol{x})$$

We can also write a neuron in terms of a dot product

$$f\left(\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_n \end{bmatrix}\right) = \sigma(\theta_0 + \boldsymbol{\theta}_{1:n} \cdot \boldsymbol{x})$$

We often write the parameters as a **weight**  $w$  and **bias**  $b$

We can also write a neuron in terms of a dot product

$$f\left(\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_n \end{bmatrix}\right) = \sigma(\theta_0 + \boldsymbol{\theta}_{1:n} \cdot \boldsymbol{x})$$

We often write the parameters as a **weight**  $\boldsymbol{w}$  and **bias**  $b$

$$f\left(\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \begin{bmatrix} b \\ w_1 \\ \vdots \\ w_n \end{bmatrix}\right) = \sigma(b + \boldsymbol{w} \cdot \boldsymbol{x})$$

We can also write a neuron in terms of a dot product

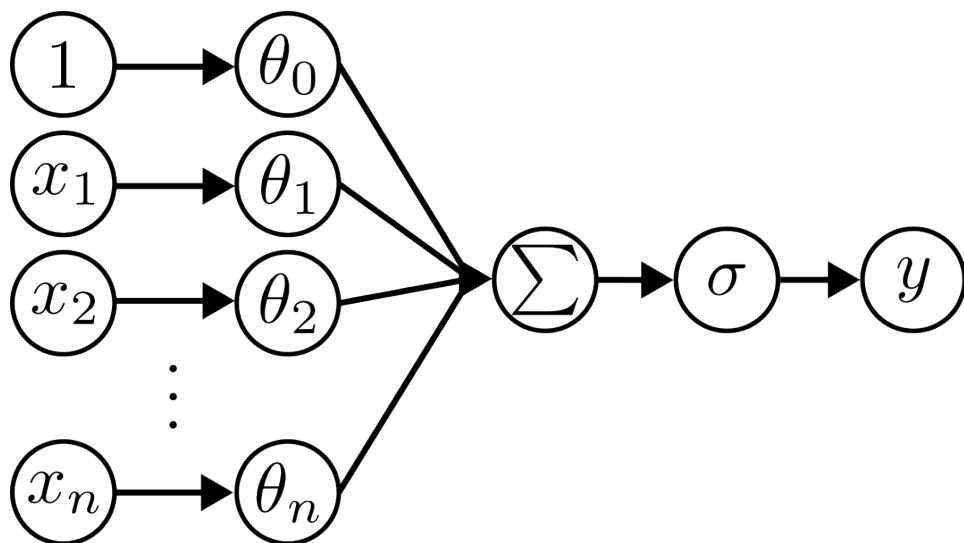
$$f\left(\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_n \end{bmatrix}\right) = \sigma(\theta_0 + \boldsymbol{\theta}_{1:n} \cdot \boldsymbol{x})$$

We often write the parameters as a **weight**  $w$  and **bias**  $b$

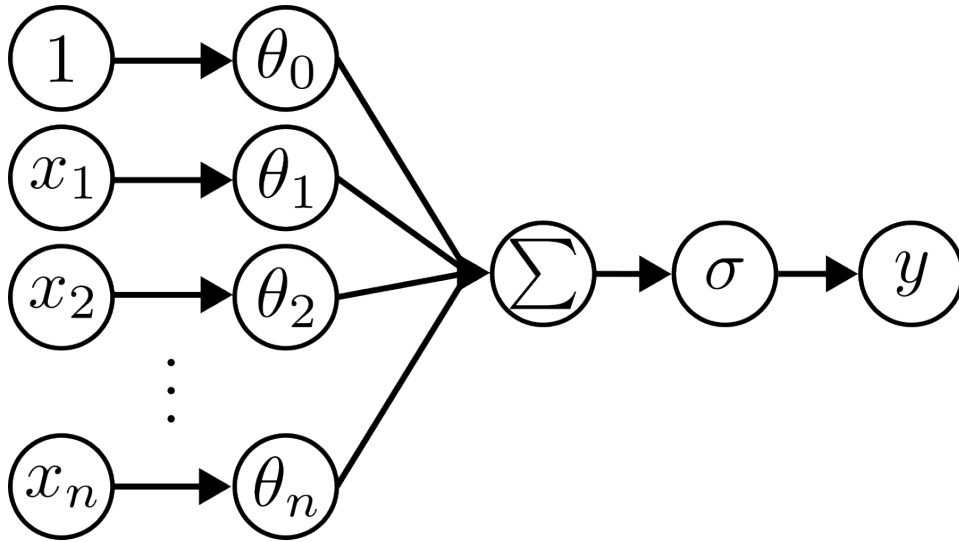
$$f\left(\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \begin{bmatrix} b \\ w_1 \\ \vdots \\ w_n \end{bmatrix}\right) = \sigma(b + \boldsymbol{w} \cdot \boldsymbol{x})$$

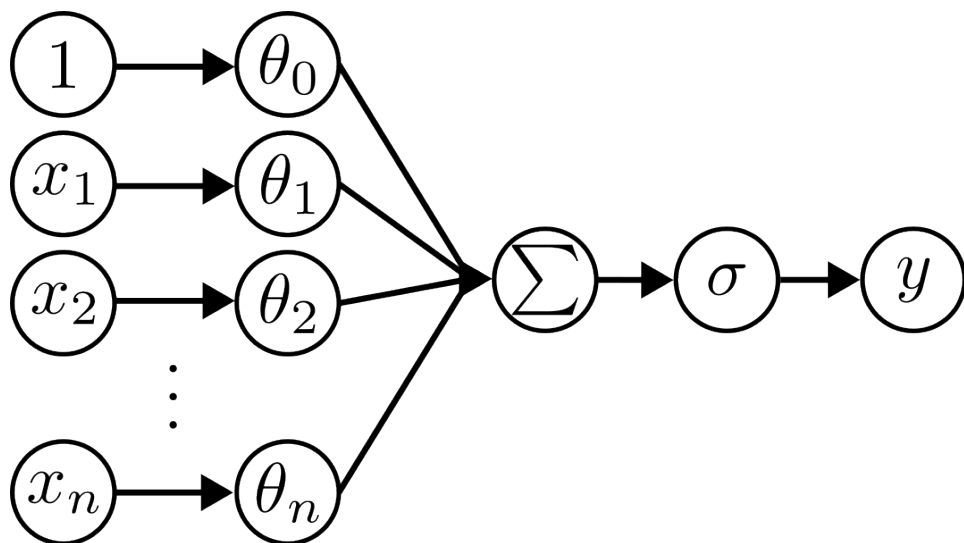
$$b = \theta_0, \boldsymbol{w} = \boldsymbol{\theta}_{1:n}$$

Relax



Recall that in machine learning we deal with functions

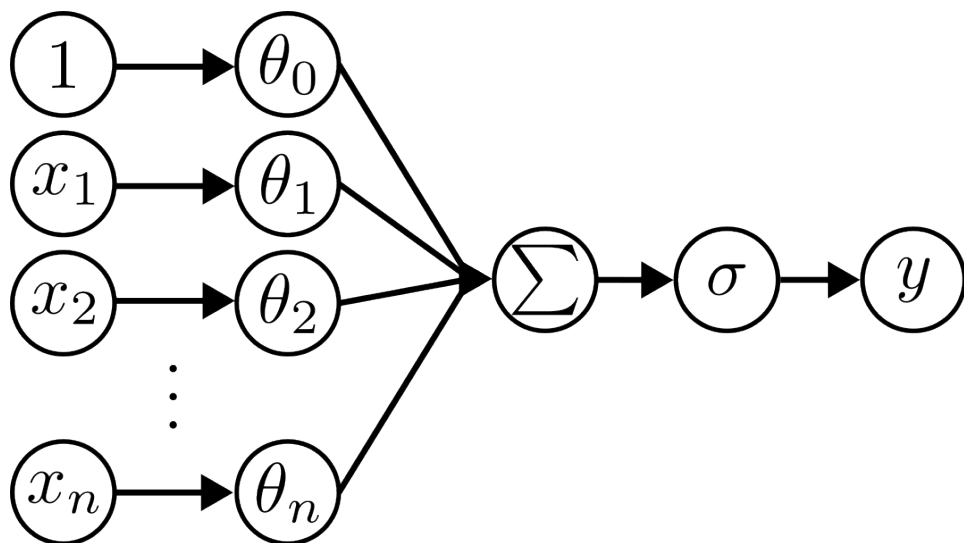




Recall that in machine learning we deal with functions

What kinds of functions can our neuron represent?



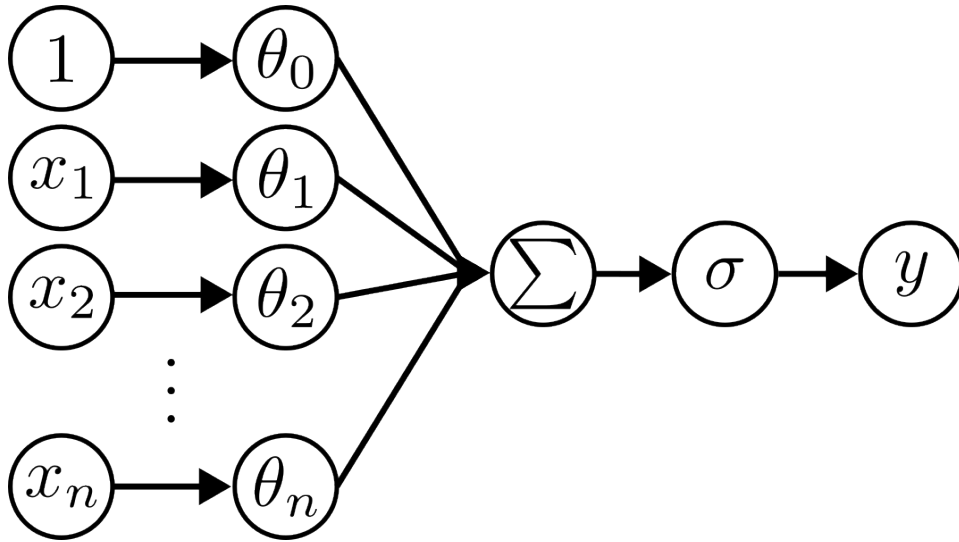


Recall that in machine learning we deal with functions

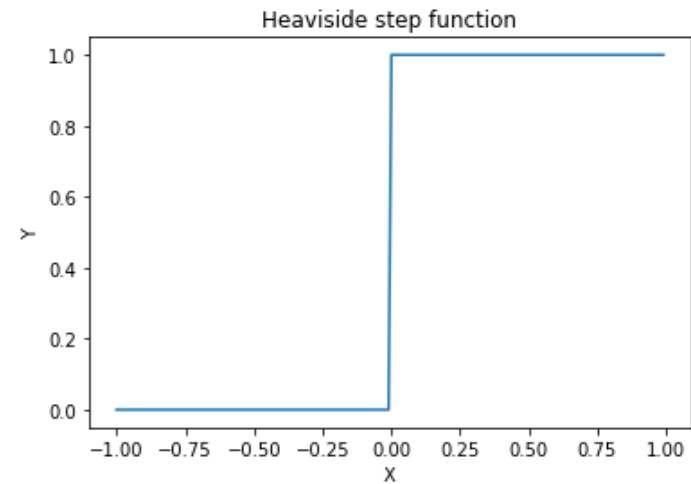
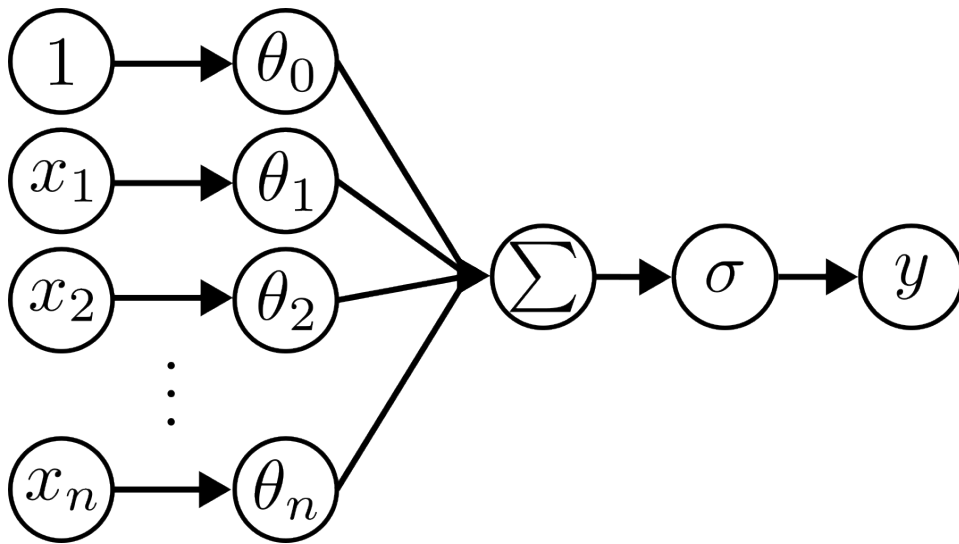
What kinds of functions can our neuron represent?

Let us start with a logical AND function

Recall the activation function  
(Heaviside step)



Recall the activation function  
(Heaviside step)



$$H(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{if } x \leq 0 \end{cases}$$

# Implement AND using an artificial neuron


Implement AND using an artificial neuron

$$f(x_1, x_2, \boldsymbol{\theta}) = H(\theta_0 + x_1\theta_1 + x_2\theta_2)$$


Implement AND using an artificial neuron

$$f(x_1, x_2, \boldsymbol{\theta}) = H(\theta_0 + x_1\theta_1 + x_2\theta_2)$$

$$\boldsymbol{\theta} = [\theta_0 \ \theta_1 \ \theta_2]^\top = [-1 \ 1 \ 1]^\top$$


## Implement AND using an artificial neuron

$$f(x_1, x_2, \boldsymbol{\theta}) = H(\theta_0 + x_1\theta_1 + x_2\theta_2)$$

$$\boldsymbol{\theta} = [\theta_0 \ \theta_1 \ \theta_2]^\top = [-1 \ 1 \ 1]^\top$$

$x_1$	$x_2$	$y$	$f(x_1, x_2, \boldsymbol{\theta})$	$\hat{y}$
0	0	0	$H(-1 + 1 \cdot 0 + 1 \cdot 0) = H(-1)$	0
0	1	0	$H(-1 + 1 \cdot 0 + 1 \cdot 1) = H(0)$	0
1	0	0	$H(-1 + 1 \cdot 1 + 1 \cdot 0) = H(0)$	0
1	1	1	$H(-1 + 1 \cdot 1 + 1 \cdot 1) = H(1)$	1

# Implement OR using an artificial neuron




## Implement OR using an artificial neuron

$$f(x_1, x_2, \boldsymbol{\theta}) = H(\theta_0 + x_1\theta_1 + x_2\theta_2)$$


## Implement OR using an artificial neuron

$$f(x_1, x_2, \boldsymbol{\theta}) = H(\theta_0 + x_1\theta_1 + x_2\theta_2)$$

$$\boldsymbol{\theta} = [\theta_0 \ \theta_1 \ \theta_2]^\top = [0 \ 1 \ 1]^\top$$


## Implement OR using an artificial neuron

$$f(x_1, x_2, \boldsymbol{\theta}) = H(\theta_0 + x_1\theta_1 + x_2\theta_2)$$

$$\boldsymbol{\theta} = [\theta_0 \ \theta_1 \ \theta_2]^\top = [0 \ 1 \ 1]^\top$$

$x_1$	$x_2$	$y$	$f(x_1, x_2, \boldsymbol{\theta})$	$\hat{y}$
0	0	0	$H(0 + 1 \cdot 0 + 1 \cdot 0) = H(0)$	0
0	1	0	$H(0 + 1 \cdot 1 + 1 \cdot 0) = H(1)$	1
1	0	1	$H(0 + 1 \cdot 0 + 1 \cdot 1) = H(1)$	1
1	1	1	$H(1 + 1 \cdot 1 + 1 \cdot 1) = H(2)$	1

# Implement XOR using an artificial neuron


## Implement XOR using an artificial neuron

$$f(x_1, x_2, \boldsymbol{\theta}) = H(\theta_0 + x_1\theta_1 + x_2\theta_2)$$


## Implement XOR using an artificial neuron

$$f(x_1, x_2, \boldsymbol{\theta}) = H(\theta_0 + x_1\theta_1 + x_2\theta_2)$$

$$\boldsymbol{\theta} = [\theta_0 \ \theta_1 \ \theta_2]^\top = [? \ ? \ ?]^\top$$


## Implement XOR using an artificial neuron

$$f(x_1, x_2, \boldsymbol{\theta}) = H(\theta_0 + x_1\theta_1 + x_2\theta_2)$$

$$\boldsymbol{\theta} = [\theta_0 \ \theta_1 \ \theta_2]^\top = [? \ ? \ ?]^\top$$

$x_1$	$x_2$	$y$	$f(x_1, x_2, \boldsymbol{\theta})$	$\hat{y}$
0	0	0	This is IMPOSSIBLE!	
0	1	1		
1	0	1		
1	1	0		

Why can't we represent XOR using a neuron?



Why can't we represent XOR using a neuron?

$$f(x_1, x_2, \boldsymbol{\theta}) = H(\theta_0 + x_1\theta_1 + x_2\theta_2)$$

Why can't we represent XOR using a neuron?

$$f(x_1, x_2, \boldsymbol{\theta}) = H(\theta_0 + x_1\theta_1 + x_2\theta_2)$$

We can only represent  $H$ (linear function)

Why can't we represent XOR using a neuron?

$$f(x_1, x_2, \boldsymbol{\theta}) = H(\theta_0 + x_1\theta_1 + x_2\theta_2)$$

We can only represent  $H$ (linear function)

XOR is not a linear combination of  $x_1, x_2$ !

Why can't we represent XOR using a neuron?

$$f(x_1, x_2, \boldsymbol{\theta}) = H(\theta_0 + x_1\theta_1 + x_2\theta_2)$$

We can only represent  $H$ (linear function)

XOR is not a linear combination of  $x_1, x_2$ !

We want to represent any function, not just linear functions

Why can't we represent XOR using a neuron?

$$f(x_1, x_2, \boldsymbol{\theta}) = H(\theta_0 + x_1\theta_1 + x_2\theta_2)$$

We can only represent  $H$ (linear function)

XOR is not a linear combination of  $x_1, x_2$ !

We want to represent any function, not just linear functions

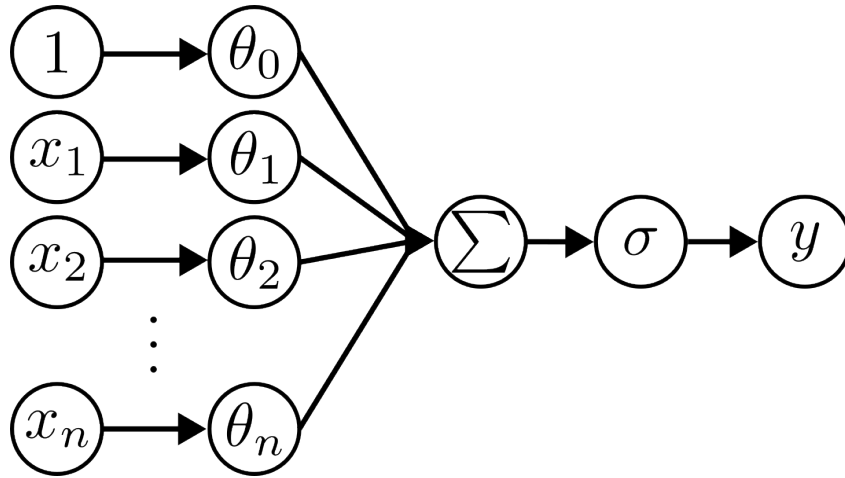
Let us think back to biology, maybe it has an answer

**Brain:** Biological neurons  $\rightarrow$  Biological neural network

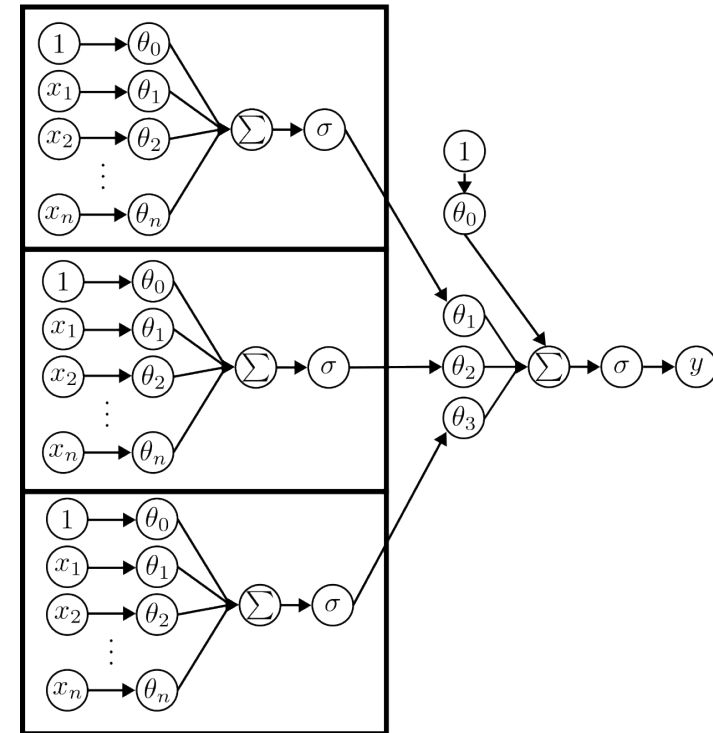
**Brain:** Biological neurons  $\rightarrow$  Biological neural network

**Computer:** Artificial neurons  $\rightarrow$  Artificial neural network

## Connect artificial neurons into a network



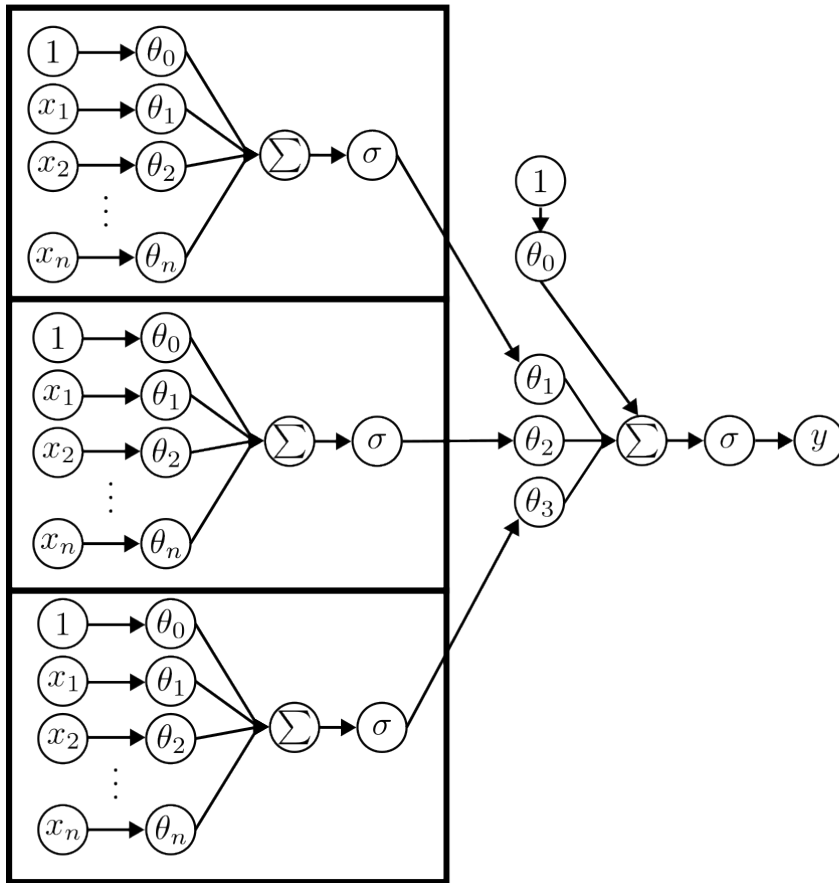
Neuron

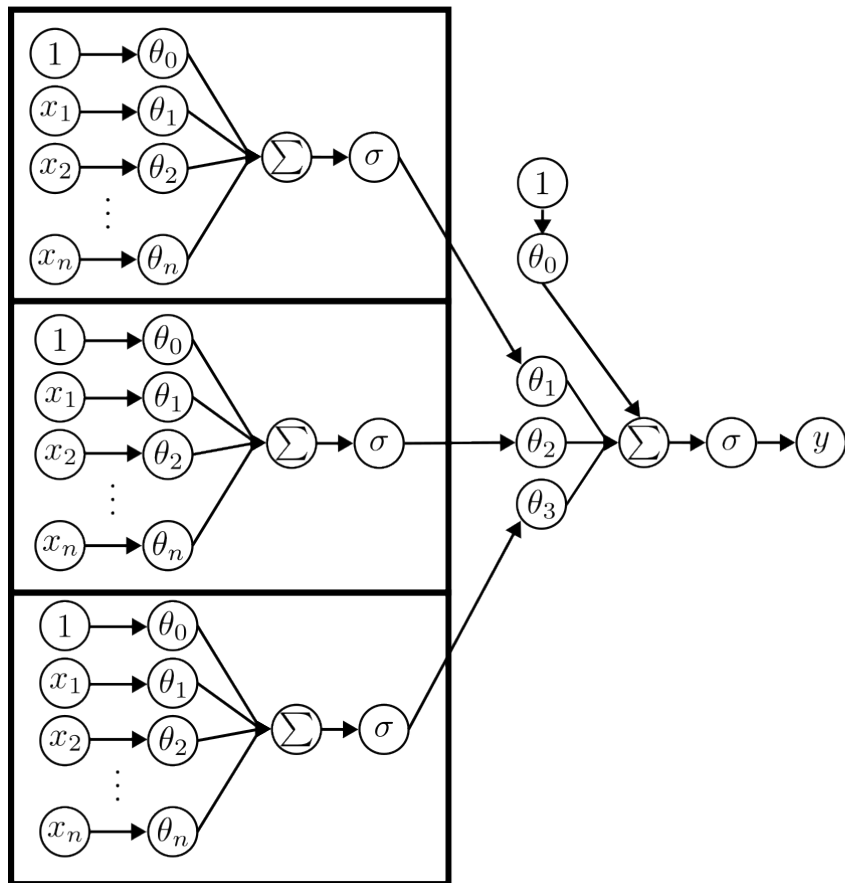


Neural Network



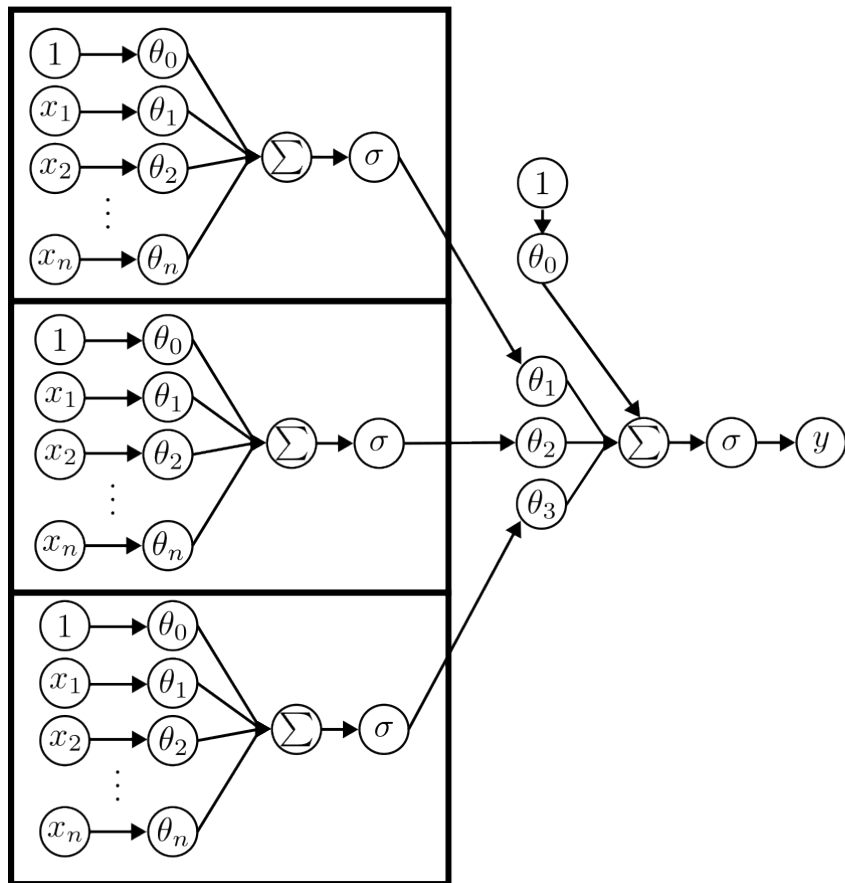
Adding neurons in **parallel**  
creates a **wide** neural network





Adding neurons in **parallel**  
creates a **wide** neural network

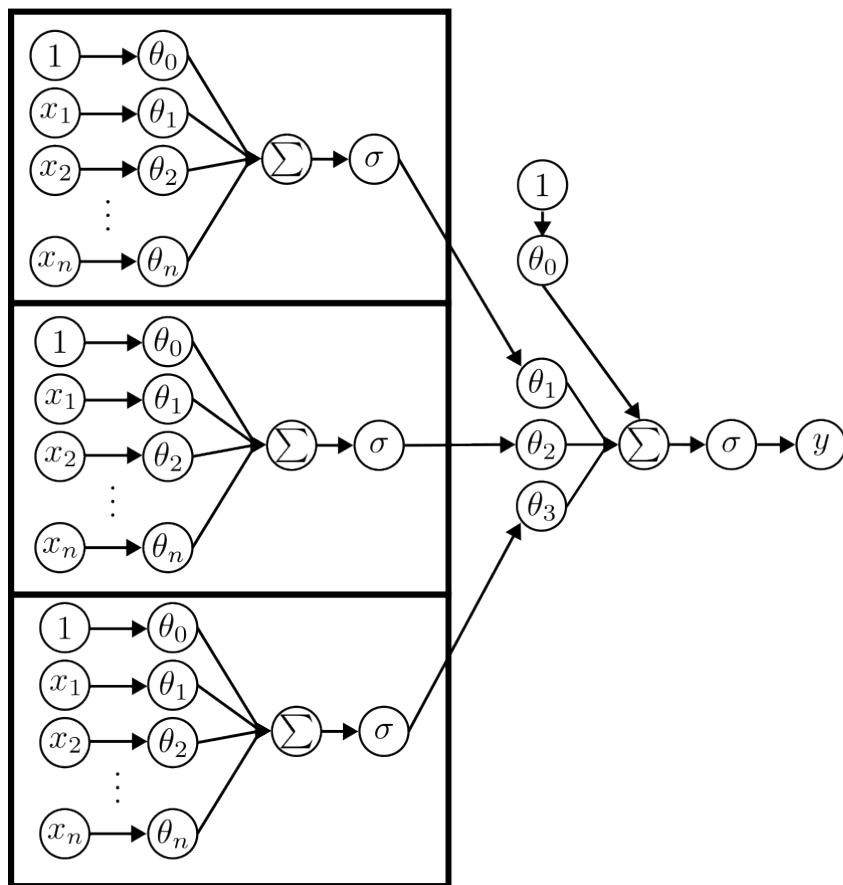
Adding neurons in **series** creates a  
**deep** neural network



Adding neurons in **parallel**  
creates a **wide** neural network

Adding neurons in **series** creates a  
**deep** neural network

Today's powerful neural networks  
are both **wide** and **deep**



Adding neurons in **parallel**  
creates a **wide** neural network

Adding neurons in **series** creates a  
**deep** neural network

Today's powerful neural networks  
are both **wide** and **deep**

Let us try to implement XOR using  
a wide and deep neural network

How do we express a **wide** neural network mathematically?

How do we express a **wide** neural network mathematically?

A single neuron:

$$f : \mathbb{R}^n, \boldsymbol{\theta} \mapsto \mathbb{R}$$

How do we express a **wide** neural network mathematically?

A single neuron:

$$f : \mathbb{R}^n, \boldsymbol{\theta} \mapsto \mathbb{R}$$

Multiple neurons (wide):

$$f : \mathbb{R}^n, \boldsymbol{\theta} \mapsto \mathbb{R}^m$$

How do we express a **wide** neural network mathematically?

A single neuron:

$$f : \mathbb{R}^n, \boldsymbol{\theta} \mapsto \mathbb{R}$$

Multiple neurons (wide):

$$f : \mathbb{R}^n, \boldsymbol{\theta} \mapsto \mathbb{R}^m$$



For a single neuron:

$$f\left(\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_n \end{bmatrix}\right) = \sigma\left(\theta_0 + \sum_{i=1}^n x_i \theta_i\right)$$

$$f\left(\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_n \end{bmatrix}\right) = \sigma(\theta_0 + \boldsymbol{\theta}_{1:n} \cdot \boldsymbol{x})$$

For a single neuron:

$$f\left(\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_n \end{bmatrix}\right) = \sigma\left(\theta_0 + \sum_{i=1}^n x_i \theta_i\right)$$

$$f\left(\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_n \end{bmatrix}\right) = \sigma(\theta_0 + \boldsymbol{\theta}_{1:n} \cdot \boldsymbol{x})$$

$$f \left( \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \begin{bmatrix} \theta_{1,0} & \theta_{2,0} & \cdots & \theta_{n,0} \\ \theta_{1,1} & \theta_{2,1} & \cdots & \theta_{n,1} \\ \vdots & \vdots & \ddots & \vdots \\ \theta_{1,m} & \theta_{2,m} & \cdots & \theta_{m,n} \end{bmatrix} \right) = \begin{bmatrix} \sigma \left( \theta_{1,0} + \sum_{i=1}^n x_i \theta_{1,i} \right) \\ \sigma \left( \theta_{2,0} + \sum_{i=1}^n x_i \theta_{2,i} \right) \\ \vdots \\ \sigma \left( \theta_{m,0} + \sum_{i=1}^n x_i \theta_{m,i} \right) \end{bmatrix}$$

Each row in the output corresponds to the output of a single neuron

$$f \left( \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \begin{bmatrix} \theta_{1,0} & \theta_{2,0} & \cdots & \theta_{n,0} \\ \theta_{1,1} & \theta_{2,1} & \cdots & \theta_{n,1} \\ \vdots & \vdots & \ddots & \vdots \\ \theta_{1,m} & \theta_{2,m} & \cdots & \theta_{m,n} \end{bmatrix} \right) = \begin{bmatrix} \sigma \left( \theta_{1,0} + \sum_{i=1}^n x_i \theta_{1,i} \right) \\ \sigma \left( \theta_{2,0} + \sum_{i=1}^n x_i \theta_{2,i} \right) \\ \vdots \\ \sigma \left( \theta_{m,0} + \sum_{i=1}^n x_i \theta_{m,i} \right) \end{bmatrix}$$

Each row in the output corresponds to the output of a single neuron

This is very confusing to write, but we can rewrite it as matrix multiplication

$$f\left(\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \begin{bmatrix} \theta_{1,0} & \theta_{2,0} & \cdots & \theta_{n,0} \\ \theta_{1,1} & \theta_{2,1} & \cdots & \theta_{n,1} \\ \vdots & \vdots & \ddots & \vdots \\ \theta_{1,m} & \theta_{2,m} & \cdots & \theta_{m,n} \end{bmatrix}\right) = \begin{bmatrix} \sigma\left(\theta_{1,0} + \sum_{i=1}^n x_i \theta_{1,i}\right) \\ \sigma\left(\theta_{2,0} + \sum_{i=1}^n x_i \theta_{2,i}\right) \\ \vdots \\ \sigma\left(\theta_{m,0} + \sum_{i=1}^n x_i \theta_{m,i}\right) \end{bmatrix}$$

$$f\left(\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \begin{bmatrix} \theta_{1,0} & \theta_{2,0} & \cdots & \theta_{n,0} \\ \theta_{1,1} & \theta_{2,1} & \cdots & \theta_{n,1} \\ \vdots & \vdots & \ddots & \vdots \\ \theta_{1,m} & \theta_{2,m} & \cdots & \theta_{m,n} \end{bmatrix}\right) = \begin{bmatrix} \sigma\left(\theta_{1,0} + \sum_{i=1}^n x_i \theta_{1,i}\right) \\ \sigma\left(\theta_{2,0} + \sum_{i=1}^n x_i \theta_{2,i}\right) \\ \vdots \\ \sigma\left(\theta_{m,0} + \sum_{i=1}^n x_i \theta_{m,i}\right) \end{bmatrix}$$

$$f(x, \theta) = \sigma(\theta_{\cdot,0} + \theta_{\cdot,1:n}x)$$

$$f\left(\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \begin{bmatrix} \theta_{1,0} & \theta_{2,0} & \cdots & \theta_{n,0} \\ \theta_{1,1} & \theta_{2,1} & \cdots & \theta_{n,1} \\ \vdots & \vdots & \ddots & \vdots \\ \theta_{1,m} & \theta_{2,m} & \cdots & \theta_{m,n} \end{bmatrix}\right) = \begin{bmatrix} \sigma\left(\theta_{1,0} + \sum_{i=1}^n x_i \theta_{1,i}\right) \\ \sigma\left(\theta_{2,0} + \sum_{i=1}^n x_i \theta_{2,i}\right) \\ \vdots \\ \sigma\left(\theta_{m,0} + \sum_{i=1}^n x_i \theta_{m,i}\right) \end{bmatrix}$$

$$f(\mathbf{x}, \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}_{\cdot,0} + \boldsymbol{\theta}_{\cdot,1:n} \mathbf{x})$$

$$f(\mathbf{x}, (\mathbf{b}, \mathbf{W})) = \sigma(\mathbf{b} + \mathbf{W}\mathbf{x})$$

How do we express a **deep** neural network mathematically?



How do we express a **deep** neural network mathematically?

A single neuron:

$$f : \mathbb{R}^n, \boldsymbol{\theta} \mapsto \mathbb{R}$$

How do we express a **deep** neural network mathematically?

A single neuron:

$$f : \mathbb{R}^n, \boldsymbol{\theta} \mapsto \mathbb{R}$$

Multiple neurons (deep):

$$f : \mathbb{R}^n, \boldsymbol{\theta}, \boldsymbol{\psi}, \dots, \boldsymbol{\rho} \mapsto \mathbb{R}^m$$

A single neuron

$$f(x, \theta) = \theta_{\cdot,0} + \theta_{\cdot,1:n}x$$

A single neuron

$$f(\mathbf{x}, \boldsymbol{\theta}) = \boldsymbol{\theta}_{\cdot,0} + \boldsymbol{\theta}_{\cdot,1:n} \mathbf{x}$$

A composition of neurons with parameters  $\boldsymbol{\theta}, \boldsymbol{\psi}, \boldsymbol{\rho}$

$$f_1(\mathbf{x}, \boldsymbol{\theta}) = \boldsymbol{\theta}_{\cdot,0} + \boldsymbol{\theta}_{\cdot,1:n} \mathbf{x} \quad f_2(\mathbf{x}, \boldsymbol{\psi}) = \boldsymbol{\psi}_{\cdot,0} + \boldsymbol{\psi}_{\cdot,1:n} \mathbf{x} \quad \dots \quad f_\ell(\mathbf{x}, \boldsymbol{\rho}) = \boldsymbol{\rho}_{\cdot,0} + \boldsymbol{\rho}_{\cdot,1:n} \mathbf{x}$$

A single neuron

$$f(\mathbf{x}, \boldsymbol{\theta}) = \boldsymbol{\theta}_{\cdot,0} + \boldsymbol{\theta}_{\cdot,1:n} \mathbf{x}$$

A composition of neurons with parameters  $\boldsymbol{\theta}, \boldsymbol{\psi}, \boldsymbol{\rho}$

$$f_1(\mathbf{x}, \boldsymbol{\theta}) = \boldsymbol{\theta}_{\cdot,0} + \boldsymbol{\theta}_{\cdot,1:n} \mathbf{x} \quad f_2(\mathbf{x}, \boldsymbol{\psi}) = \boldsymbol{\psi}_{\cdot,0} + \boldsymbol{\psi}_{\cdot,1:n} \mathbf{x} \quad \dots \quad f_\ell(\mathbf{x}, \boldsymbol{\rho}) = \boldsymbol{\rho}_{\cdot,0} + \boldsymbol{\rho}_{\cdot,1:n} \mathbf{x}$$

A single neuron

$$f(\mathbf{x}, \boldsymbol{\theta}) = \boldsymbol{\theta}_{\cdot,0} + \boldsymbol{\theta}_{\cdot,1:n} \mathbf{x}$$

A composition of neurons with parameters  $\boldsymbol{\theta}, \boldsymbol{\psi}, \boldsymbol{\rho}$

$$f_1(\mathbf{x}, \boldsymbol{\theta}) = \boldsymbol{\theta}_{\cdot,0} + \boldsymbol{\theta}_{\cdot,1:n} \mathbf{x} \quad f_2(\mathbf{x}, \boldsymbol{\psi}) = \boldsymbol{\psi}_{\cdot,0} + \boldsymbol{\psi}_{\cdot,1:n} \mathbf{x} \quad \dots \quad f_\ell(\mathbf{x}, \boldsymbol{\rho}) = \boldsymbol{\rho}_{\cdot,0} + \boldsymbol{\rho}_{\cdot,1:n} \mathbf{x}$$

$$f_\ell(\dots f_2(f_1(\mathbf{x}, \boldsymbol{\theta}_1), \boldsymbol{\psi}) \dots)$$

Written more plainly as

$$z_1 = f_1(x, \theta) = \theta_{\cdot,0} + \theta_{\cdot,1:n}x$$

Written more plainly as

$$z_1 = f_1(x, \theta) = \theta_{\cdot,0} + \theta_{\cdot,1:n}x$$

$$z_2 = f_2(z_1, \psi) = \psi_{\cdot,0} + \psi_{\cdot,1:n}z_1$$



Written more plainly as

$$z_1 = f_1(x, \theta) = \theta_{\cdot,0} + \theta_{\cdot,1:n}x$$

$$z_2 = f_2(z_1, \psi) = \psi_{\cdot,0} + \psi_{\cdot,1:n}z_1$$

$$\vdots$$

Written more plainly as

$$z_1 = f_1(x, \theta) = \theta_{\cdot,0} + \theta_{\cdot,1:n}x$$

$$z_2 = f_2(z_1, \psi) = \psi_{\cdot,0} + \psi_{\cdot,1:n}z_1$$

$$\vdots$$

$$y = f_\ell(x, \rho) = \rho_{\cdot,0} + \rho_{\cdot,1:n}z_{\ell-1}$$

Written more plainly as

$$z_1 = f_1(x, \theta) = \theta_{\cdot,0} + \theta_{\cdot,1:n}x$$

$$z_2 = f_2(z_1, \psi) = \psi_{\cdot,0} + \psi_{\cdot,1:n}z_1$$

$$\vdots$$

$$y = f_\ell(x, \rho) = \rho_{\cdot,0} + \rho_{\cdot,1:n}z_{\ell-1}$$

# Implement XOR using a deep and wide neural network

Implement XOR using a deep and wide neural network

$$f(x_1, x_2, \boldsymbol{\theta}) = H(\theta_{3,0} + \theta_{3,1} \cdot H(\theta_{1,0} + x_1\theta_{1,1} + x_2\theta_{1,2}) + \theta_{3,2} \cdot H(\theta_{2,0} + x_1\theta_{2,1} + x_2\theta_{2,2}))$$

Implement XOR using a deep and wide neural network

$$f(x_1, x_2, \boldsymbol{\theta}) = H(\theta_{3,0} + \theta_{3,1} \cdot H(\theta_{1,0} + x_1\theta_{1,1} + x_2\theta_{1,2}) + \theta_{3,2} \cdot H(\theta_{2,0} + x_1\theta_{2,1} + x_2\theta_{2,2}))$$

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_{1,0} & \theta_{1,1} & \theta_{1,2} \\ \theta_{2,0} & \theta_{2,1} & \theta_{2,2} \\ \theta_{3,0} & \theta_{3,1} & \theta_{3,2} \end{bmatrix} = \begin{bmatrix} -0.5 & 1 & 1 \\ -1.5 & 1 & 1 \\ -0.5 & 1 & -2 \end{bmatrix}$$

Implement XOR using a deep and wide neural network

$$f(x_1, x_2, \boldsymbol{\theta}) = H(\theta_{3,0} + \theta_{3,1} \cdot H(\theta_{1,0} + x_1\theta_{1,1} + x_2\theta_{1,2}) + \theta_{3,2} \cdot H(\theta_{2,0} + x_1\theta_{2,1} + x_2\theta_{2,2}))$$

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_{1,0} & \theta_{1,1} & \theta_{1,2} \\ \theta_{2,0} & \theta_{2,1} & \theta_{2,2} \\ \theta_{3,0} & \theta_{3,1} & \theta_{3,2} \end{bmatrix} = \begin{bmatrix} -0.5 & 1 & 1 \\ -1.5 & 1 & 1 \\ -0.5 & 1 & -2 \end{bmatrix}$$

What other functions can we represent using a deep and wide neural network?



What other functions can we represent using a deep and wide neural network?

Consider a one-dimensional arbitrary function  $g(x) = y$

What other functions can we represent using a deep and wide neural network?

Consider a one-dimensional arbitrary function  $g(x) = y$

We can approximate  $g$  using our neural network  $f$

What other functions can we represent using a deep and wide neural network?

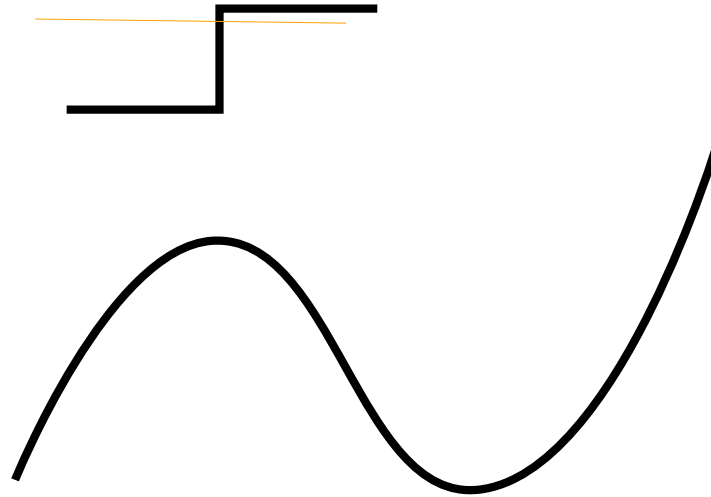
Consider a one-dimensional arbitrary function  $g(x) = y$

We can approximate  $g$  using our neural network  $f$

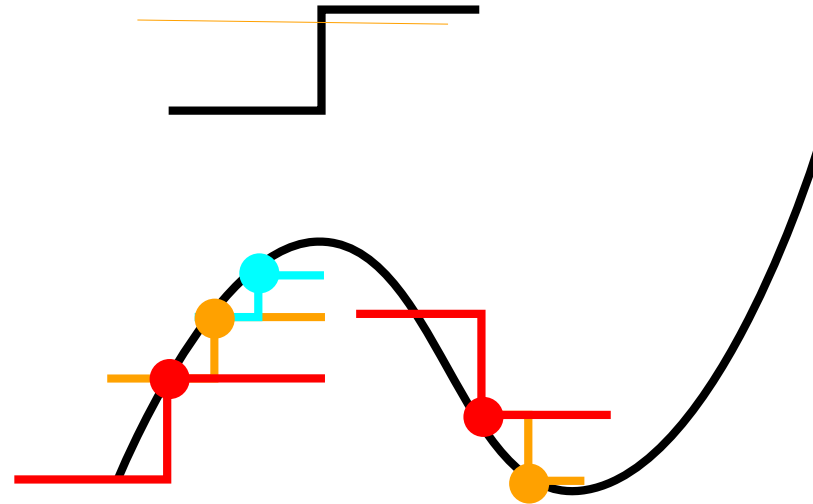
$$\begin{aligned} f(x_1, x_2, \boldsymbol{\theta}) = & H(\theta_{3,0} \\ & + \theta_{3,1} \cdot H(\theta_{1,0} + x_1 \theta_{1,1} + x_2 \theta_{1,2}) \\ & + \theta_{3,2} \cdot H(\theta_{2,0} + x_1 \theta_{2,1} + x_2 \theta_{2,2})) \end{aligned}$$

**Proof Sketch:** Approximate a function  $g(x)$  using a linear combination of Heaviside functions

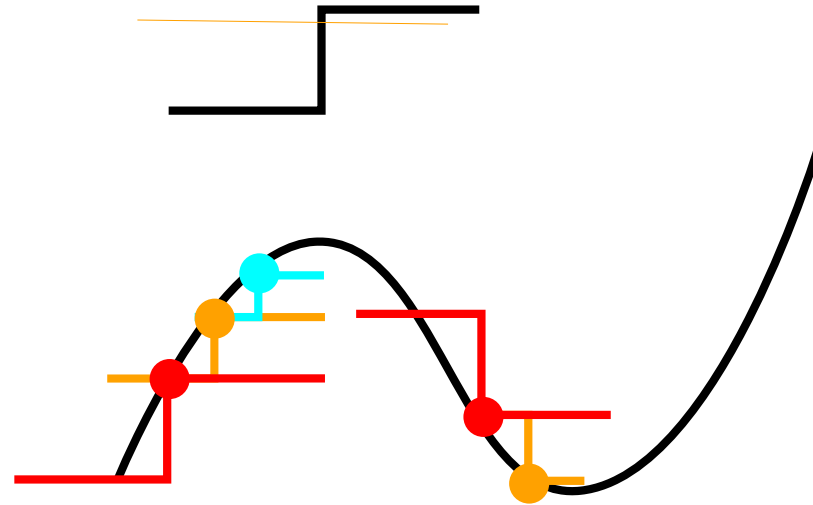
**Proof Sketch:** Approximate a function  $g(x)$  using a linear combination of Heaviside functions



**Proof Sketch:** Approximate a function  $g(x)$  using a linear combination of Heaviside functions



**Proof Sketch:** Approximate a function  $g(x)$  using a linear combination of Heaviside functions



$$\text{Roughly, } \exists \boldsymbol{\theta} \Rightarrow \lim_{n \rightarrow \infty} \left[ \theta_{2,0} + \theta_{2,1} \sum_{j=1}^n \sigma(\theta_{1,0} + \theta_{1,j}x) \right] = g(x); \quad \forall g$$

More formally, a wide and deep neural network is a **universal function approximator**



More formally, a wide and deep neural network is a **universal function approximator**

It can approximate **any** continuous function to precision  $\varepsilon$

More formally, a wide and deep neural network is a **universal function approximator**

It can approximate **any** continuous function to precision  $\varepsilon$

$$| g(\boldsymbol{x}) - f(\boldsymbol{x}, \boldsymbol{\theta}) | < \varepsilon$$

More formally, a wide and deep neural network is a **universal function approximator**

It can approximate **any** continuous function to precision  $\varepsilon$

$$| g(\boldsymbol{x}) - f(\boldsymbol{x}, \boldsymbol{\theta}) | < \varepsilon$$

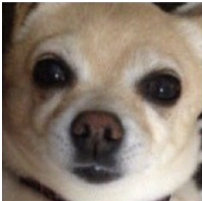
As we increase the width and depth of the network,  $\varepsilon$  shrinks

More formally, a wide and deep neural network is a **universal function approximator**

It can approximate **any** continuous function to precision  $\varepsilon$

$$| g(\mathbf{x}) - f(\mathbf{x}, \boldsymbol{\theta}) | < \varepsilon$$

As we increase the width and depth of the network,  $\varepsilon$  shrinks

$$g\left(\text{$$
 = Dog
$$g\left(\text{$$
 = Muffin

Very powerful finding! The basis of deep learning.

All the models we examine in this course will use this neural network structure internally

All the models we examine in this course will use this neural network structure internally

- Transformers

All the models we examine in this course will use this neural network structure internally

- Transformers
- Graph neural networks

All the models we examine in this course will use this neural network structure internally

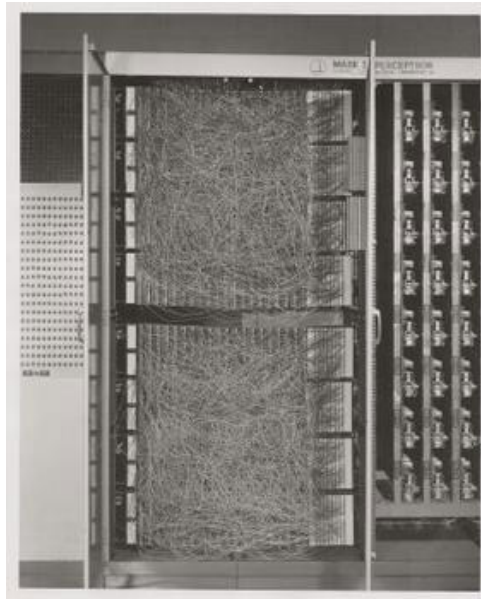
- Transformers
- Graph neural networks



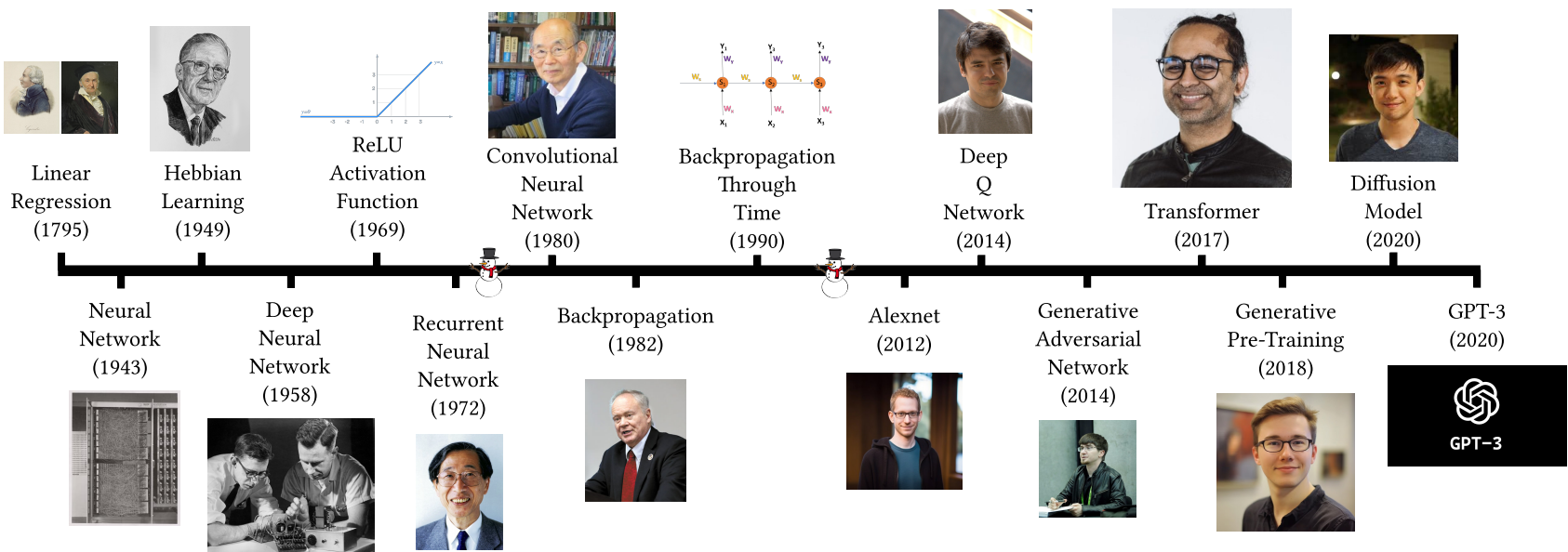
Relax

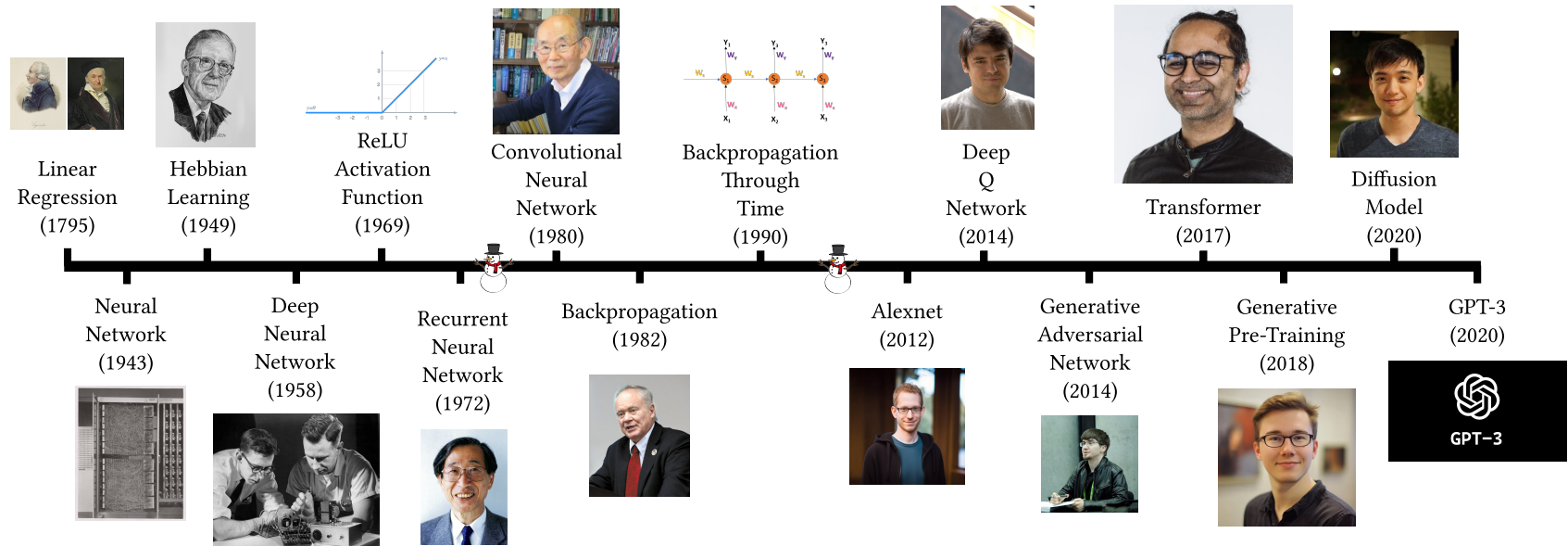
We call this form of a neural network a **feedforward network** or **perceptron** (invented in 1943)

We call this form of a neural network a **feedforward network** or **perceptron** (invented in 1943)



$20 \times 20$  grid of pixels to process images





**Question:** If the deep neural network was invented in 1958, why did it take 70 years for us to care about deep learning?

**Answer:** Deep learning requires very deep and wide networks

1. Hardware advances enabled very deep and wide networks

**Answer:** Deep learning requires very deep and wide networks

1. Hardware advances enabled very deep and wide networks
2. Many theoretical improvements allow us to successfully train deeper and wider networks

**Answer:** Deep learning requires very deep and wide networks

1. Hardware advances enabled very deep and wide networks
2. Many theoretical improvements allow us to successfully train deeper and wider networks



The neural network we created today is called a feedforward network or perceptron

The neural network we created today is called a feedforward network or perceptron

When the network is deep, we call it a Multi-Layer Perceptron (MLP)

The neural network we created today is called a feedforward network or perceptron

When the network is deep, we call it a Multi-Layer Perceptron (MLP)

We often use the term “layers”, when referring to a specific depth of the neural network

- Four-layer MLP means a neural network with a depth of four
- Corresponds to four parameter matrices in  $\theta$

Let us construct deep and wide neural networks in torch and jax

Here are the equations for one neural network layer

$$f(\mathbf{x}, \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}_{\cdot,0} + \boldsymbol{\theta}_{\cdot,1:n} \mathbf{x}) \quad \text{or} \quad f(\mathbf{x}, (\mathbf{b}, \mathbf{W})) = \sigma(\mathbf{b} + \mathbf{W} \mathbf{x})$$

Here are the equations for one neural network layer

$$f(\mathbf{x}, \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}_{:,0} + \boldsymbol{\theta}_{:,1:n} \mathbf{x}) \quad \text{or} \quad f(\mathbf{x}, (\mathbf{b}, \mathbf{W})) = \sigma(\mathbf{b} + \mathbf{W} \mathbf{x})$$

We must implement the linear function  $\mathbf{b} + \mathbf{W} \mathbf{x}$  and the activation function  $\sigma$  to create a neural network layer

Here are the equations for one neural network layer

$$f(\mathbf{x}, \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}_{:,0} + \boldsymbol{\theta}_{:,1:n} \mathbf{x}) \quad \text{or} \quad f(\mathbf{x}, (\mathbf{b}, \mathbf{W})) = \sigma(\mathbf{b} + \mathbf{W} \mathbf{x})$$

We must implement the linear function  $\mathbf{b} + \mathbf{W} \mathbf{x}$  and the activation function  $\sigma$  to create a neural network layer

Let us do this in colab! [https://colab.research.google.com/drive/1bLtf3QY-yROlif\\_EoQSU1WS7svd0q8j7?usp=sharing](https://colab.research.google.com/drive/1bLtf3QY-yROlif_EoQSU1WS7svd0q8j7?usp=sharing)