

Neural Networks

CISC 7026: Introduction to Deep Learning

University of Macau

1. Review
2. Multivariate linear regression
3. Limitations of linear regression
4. History of neural networks
5. Biological neurons
6. Artificial neurons
7. Perceptron
8. Multilayer Perceptron

1. Review
2. Multivariate linear regression
3. Limitations of linear regression
4. History of neural networks
5. Biological neurons
6. Artificial neurons
7. Perceptron
8. Multilayer Perceptron

Review

Since you are very educated, we focused on how education affects life expectancy

Review

Since you are very educated, we focused on how education affects life expectancy

Studies show a causal effect of education on health

Since you are very educated, we focused on how education affects life expectancy

Studies show a causal effect of education on health

- *The causal effects of education on health outcomes in the UK Biobank.*
Davies et al. *Nature Human Behaviour.*

Since you are very educated, we focused on how education affects life expectancy

Studies show a causal effect of education on health

- *The causal effects of education on health outcomes in the UK Biobank.*
Davies et al. *Nature Human Behaviour*.
- By staying in school, you are likely to live longer

Since you are very educated, we focused on how education affects life expectancy

Studies show a causal effect of education on health

- *The causal effects of education on health outcomes in the UK Biobank.*
Davies et al. *Nature Human Behaviour*.
- By staying in school, you are likely to live longer
- Being rich also helps, but education alone has a **causal** relationship with life expectancy

Task: Given your education, predict your life expectancy

Task: Given your education, predict your life expectancy

$X \in \mathbb{R}_+$: Years in school

Task: Given your education, predict your life expectancy

$X \in \mathbb{R}_+$: Years in school

$Y \in \mathbb{R}_+$: Age of death

Task: Given your education, predict your life expectancy

$X \in \mathbb{R}_+$: Years in school

$Y \in \mathbb{R}_+$: Age of death

$\Theta \in \mathbb{R}^2$: Parameters

Task: Given your education, predict your life expectancy

$X \in \mathbb{R}_+$: Years in school

$Y \in \mathbb{R}_+$: Age of death

$\Theta \in \mathbb{R}^2$: Parameters

$$f : X \times \Theta \mapsto Y$$

Approach: Learn the parameters θ such that

$$f(x, \theta) = y; \quad x \in X, y \in Y$$

Task: Given your education, predict your life expectancy

$X \in \mathbb{R}_+$: Years in school

$Y \in \mathbb{R}_+$: Age of death

$\Theta \in \mathbb{R}^2$: Parameters

$$f : X \times \Theta \mapsto Y$$

Approach: Learn the parameters θ such that

$$f(x, \theta) = y; \quad x \in X, y \in Y$$

Goal: Given someone's education, predict how long they will live

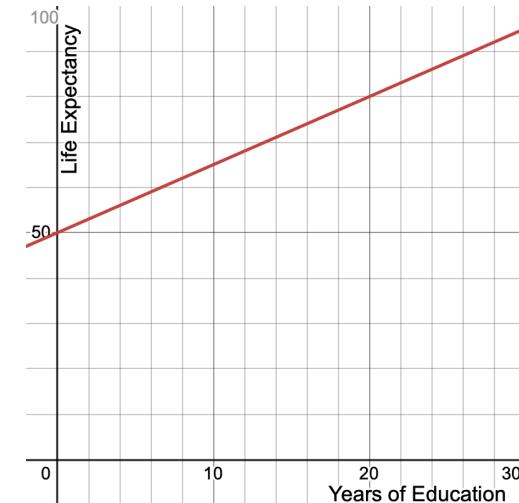
Review

Started with a linear function f

Review

Started with a linear function f

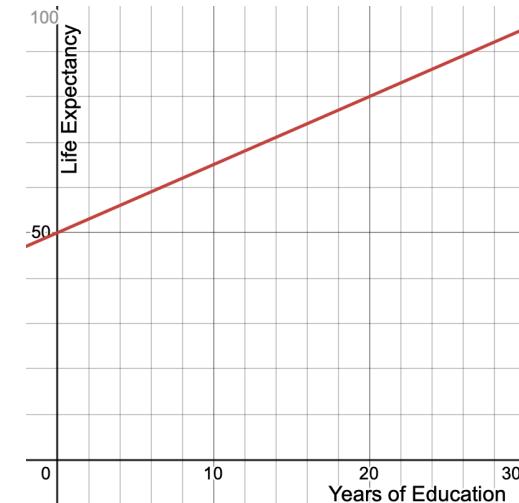
$$f(x, \theta) = f\left(x, \begin{bmatrix} \theta_1 \\ \theta_0 \end{bmatrix}\right) = \theta_1 x + \theta_0$$



Review

Started with a linear function f

$$f(x, \theta) = f\left(x, \begin{bmatrix} \theta_1 \\ \theta_0 \end{bmatrix}\right) = \theta_1 x + \theta_0$$

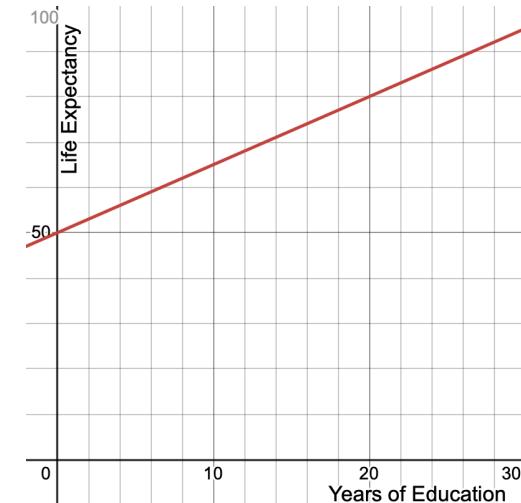


Then, we derived the square error function

Review

Started with a linear function f

$$f(x, \theta) = f\left(x, \begin{bmatrix} \theta_1 \\ \theta_0 \end{bmatrix}\right) = \theta_1 x + \theta_0$$



Then, we derived the square error function

$$\text{error}(f(x, \theta), y) = (f(x, \theta) - y)^2$$

Review

We wrote the loss function for a single datapoint x_i, y_i using the square error

$$\mathcal{L}(x_i, y_i, \theta) = \text{error}(f(x_i, \theta), y_i) = (f(x_i, \theta) - y_i)^2$$

Review

We wrote the loss function for a single datapoint x_i, y_i using the square error

$$\mathcal{L}(x_i, y_i, \theta) = \text{error}(f(x_i, \theta), y_i) = (f(x_i, \theta) - y_i)^2$$

But we wanted to learn a model over **all** the data, not a single datapoint

Review

We wrote the loss function for a single datapoint x_i, y_i using the square error

$$\mathcal{L}(x_i, y_i, \theta) = \text{error}(f(x_i, \theta), y_i) = (f(x_i, \theta) - y_i)^2$$

But we wanted to learn a model over **all** the data, not a single datapoint

We wanted to make **new** predictions, to **generalize**

Review

We wrote the loss function for a single datapoint x_i, y_i using the square error

$$\mathcal{L}(x_i, y_i, \theta) = \text{error}(f(x_i, \theta), y_i) = (f(x_i, \theta) - y_i)^2$$

But we wanted to learn a model over **all** the data, not a single datapoint

We wanted to make **new** predictions, to **generalize**

$$\mathbf{x} = [x_1 \ x_2 \ \dots \ x_n]^\top, \mathbf{y} = [y_1 \ y_2 \ \dots \ y_n]^\top$$

Review

We wrote the loss function for a single datapoint x_i, y_i using the square error

$$\mathcal{L}(x_i, y_i, \theta) = \text{error}(f(x_i, \theta), y_i) = (f(x_i, \theta) - y_i)^2$$

But we wanted to learn a model over **all** the data, not a single datapoint

We wanted to make **new** predictions, to **generalize**

$$\mathbf{x} = [x_1 \ x_2 \ \dots \ x_n]^\top, \mathbf{y} = [y_1 \ y_2 \ \dots \ y_n]^\top$$

$$\mathcal{L}(\mathbf{x}, \mathbf{y}, \theta) = \sum_{i=1}^n \text{error}(f(x_i, \theta), y_i) = \sum_{i=1}^n (f(x_i, \theta) - y_i)^2$$

Review

Our objective was to find the parameters that minimized the loss function over the dataset

Review

Our objective was to find the parameters that minimized the loss function over the dataset

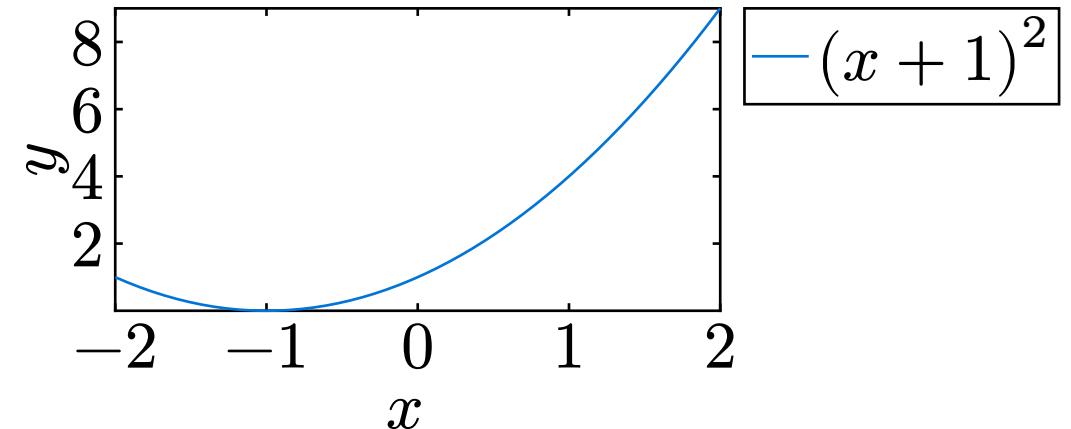
We introduced the arg min operator

Review

Our objective was to find the parameters that minimized the loss function over the dataset

We introduced the arg min operator

$$f(x) = (x + 1)^2$$

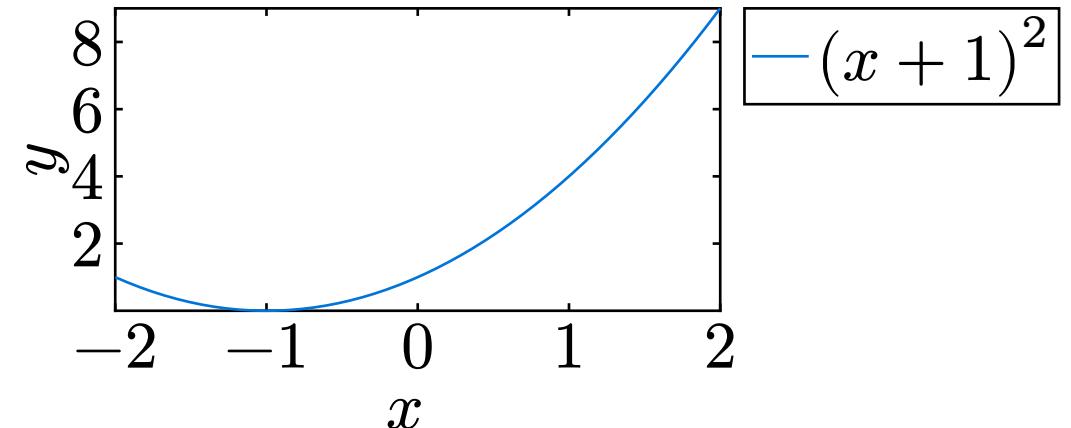


Review

Our objective was to find the parameters that minimized the loss function over the dataset

We introduced the arg min operator

$$f(x) = (x + 1)^2$$



$$\arg \min_x f(x) = -1$$

With the $\arg \min$ operator, we formally wrote our optimization objective

With the $\arg \min$ operator, we formally wrote our optimization objective

$$\begin{aligned}\arg \min_{\theta} \mathcal{L}(x, y, \theta) &= \arg \min_{\theta} \sum_{i=1}^n \text{error}(f(x_i, \theta), y_i) \\ &= \arg \min_{\theta} \sum_{i=1}^n (f(x_i, \theta) - y_i)^2\end{aligned}$$

With the $\arg \min$ operator, we formally wrote our optimization objective

$$\begin{aligned}\arg \min_{\theta} \mathcal{L}(x, y, \theta) &= \arg \min_{\theta} \sum_{i=1}^n \text{error}(f(x_i, \theta), y_i) \\ &= \arg \min_{\theta} \sum_{i=1}^n (f(x_i, \theta) - y_i)^2\end{aligned}$$

Review

We defined the design matrix X_D

Review

We defined the design matrix \mathbf{X}_D

$$\mathbf{X}_D = [x \ 1] = \begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ \vdots & \vdots \\ x_n & 1 \end{bmatrix}$$

We defined the design matrix \mathbf{X}_D

$$\mathbf{X}_D = [x \ 1] = \begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ \vdots & \vdots \\ x_n & 1 \end{bmatrix}$$

With the design matrix, provided an **analytical** solution to the optimization objective

We defined the design matrix \mathbf{X}_D

$$\mathbf{X}_D = [\mathbf{x} \ 1] = \begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ \vdots & \vdots \\ x_n & 1 \end{bmatrix}$$

With the design matrix, provided an **analytical** solution to the optimization objective

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

We defined the design matrix \mathbf{X}_D

$$\mathbf{X}_D = [\mathbf{x} \ 1] = \begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ \vdots & \vdots \\ x_n & 1 \end{bmatrix}$$

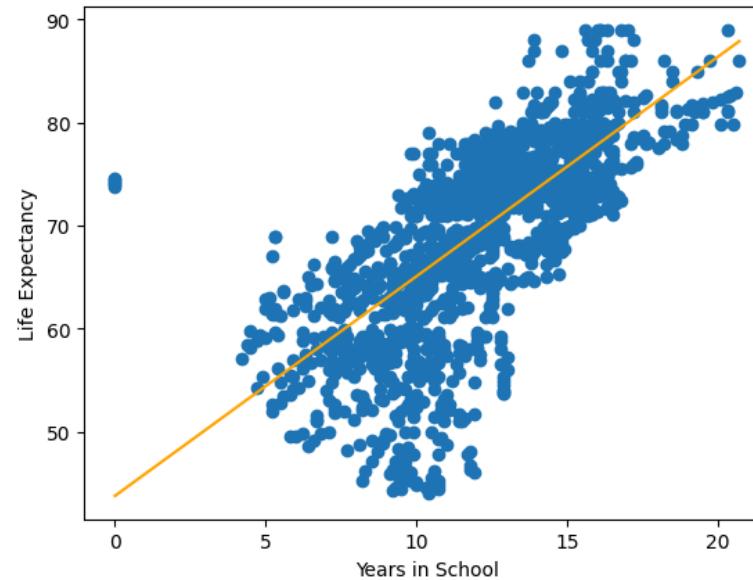
With the design matrix, provided an **analytical** solution to the optimization objective

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

With this analytical solution, we were able to learn a linear model

Review

With this analytical solution, we were able to learn a linear model



Then, we used a trick to extend linear regression to nonlinear models

Then, we used a trick to extend linear regression to nonlinear models

$$\mathbf{X}_D = \begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ \vdots & \vdots \\ x_n & 1 \end{bmatrix} \Rightarrow \mathbf{X}_D = \begin{bmatrix} \log(1 + x_1) & 1 \\ \log(1 + x_2) & 1 \\ \vdots & \vdots \\ \log(1 + x_n) & 1 \end{bmatrix}$$

Then, we used a trick to extend linear regression to nonlinear models

$$\mathbf{X}_D = \begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ \vdots & \vdots \\ x_n & 1 \end{bmatrix} \Rightarrow \mathbf{X}_D = \begin{bmatrix} \log(1 + x_1) & 1 \\ \log(1 + x_2) & 1 \\ \vdots & \vdots \\ \log(1 + x_n) & 1 \end{bmatrix}$$

We extended to polynomial, which are **universal function approximators**

We extended to polynomial, which are **universal function approximators**

$$\mathbf{X}_D = \begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ \vdots & \vdots \\ x_n & 1 \end{bmatrix} \Rightarrow \mathbf{X}_D = \begin{bmatrix} x_1^m & x_1^{m-1} & \dots & x_1 & 1 \\ x_2^m & x_2^{m-1} & \dots & x_2 & 1 \\ \vdots & \vdots & \ddots & & \\ x_n^m & x_n^{m-1} & \dots & x_n & 1 \end{bmatrix}$$

We extended to polynomial, which are **universal function approximators**

$$\mathbf{X}_D = \begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ \vdots & \vdots \\ x_n & 1 \end{bmatrix} \Rightarrow \mathbf{X}_D = \begin{bmatrix} x_1^m & x_1^{m-1} & \dots & x_1 & 1 \\ x_2^m & x_2^{m-1} & \dots & x_2 & 1 \\ \vdots & \vdots & \ddots & & \\ x_n^m & x_n^{m-1} & \dots & x_n & 1 \end{bmatrix}$$

$$f : X \times \Theta \mapsto \mathbb{R}$$

We extended to polynomial, which are **universal function approximators**

$$\mathbf{X}_D = \begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ \vdots & \vdots \\ x_n & 1 \end{bmatrix} \Rightarrow \mathbf{X}_D = \begin{bmatrix} x_1^m & x_1^{m-1} & \dots & x_1 & 1 \\ x_2^m & x_2^{m-1} & \dots & x_2 & 1 \\ \vdots & \vdots & \ddots & & \\ x_n^m & x_n^{m-1} & \dots & x_n & 1 \end{bmatrix}$$

$$f : X \times \Theta \mapsto \mathbb{R}$$

$$\Theta \in \mathbb{R}^2 \Rightarrow \Theta \in \mathbb{R}^m$$

Finally, we discussed overfitting

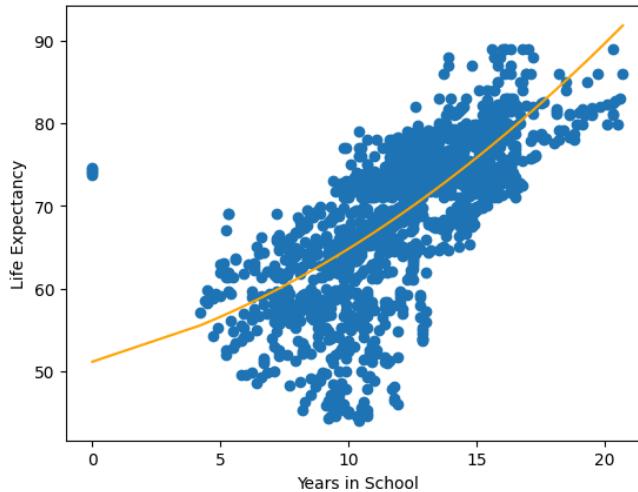
Finally, we discussed overfitting

$$f(x, \theta) = \theta_n x^m + \theta_{m-1} x^{m-1}, \dots, \theta_1 + x^1 + \theta_0$$

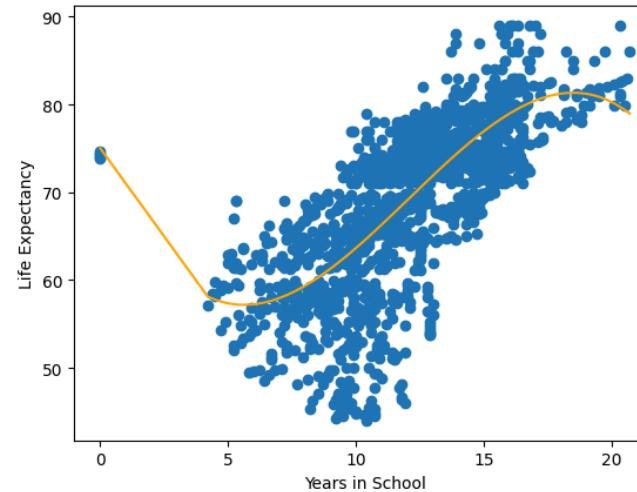
Review

Finally, we discussed overfitting

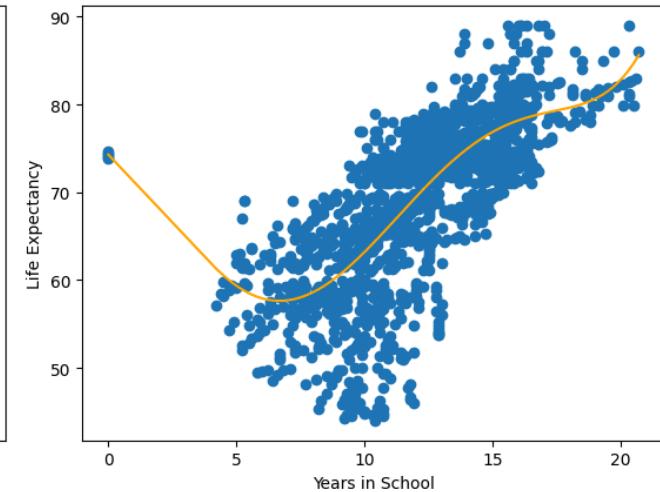
$$f(x, \theta) = \theta_n x^m + \theta_{m-1} x^{m-1}, \dots, \theta_1 + x^1 + \theta_0$$



$$m = 2$$



$$m = 3$$



$$m = 5$$

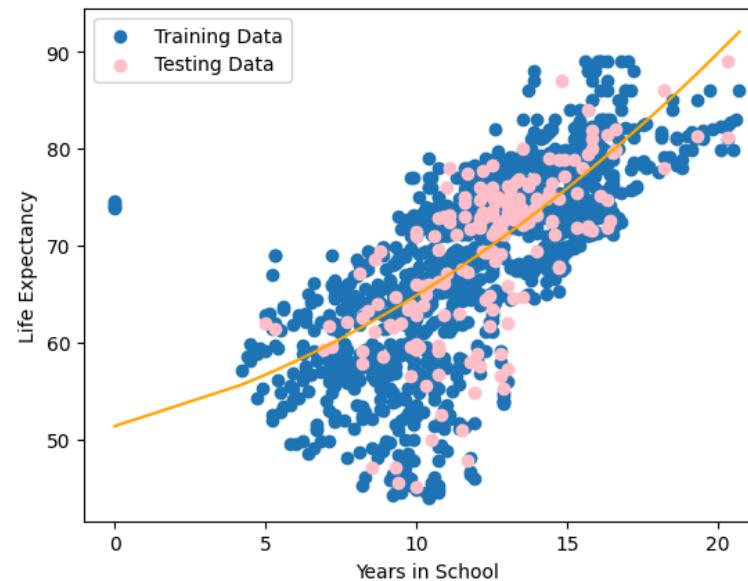
We care about **generalization** in machine learning

We care about **generalization** in machine learning

So we should always split our dataset into a training dataset and a testing dataset

We care about **generalization** in machine learning

So we should always split our dataset into a training dataset and a testing dataset



1. Review
2. Multivariate linear regression
3. Limitations of linear regression
4. History of neural networks
5. Biological neurons
6. Artificial neurons
7. Perceptron
8. Multilayer Perceptron

1. Review
2. **Multivariate linear regression**
3. Limitations of linear regression
4. History of neural networks
5. Biological neurons
6. Artificial neurons
7. Perceptron
8. Multilayer Perceptron

Last time, we assumed a single-input system

Last time, we assumed a single-input system

Years of education: $X \in \mathbb{R}$

Last time, we assumed a single-input system

Years of education: $X \in \mathbb{R}$

But sometimes we want to consider multiple input dimensions

Last time, we assumed a single-input system

Years of education: $X \in \mathbb{R}$

But sometimes we want to consider multiple input dimensions

Years of education, BMI, GDP: $X \in \mathbb{R}^3$

Last time, we assumed a single-input system

Years of education: $X \in \mathbb{R}$

But sometimes we want to consider multiple input dimensions

Years of education, BMI, GDP: $X \in \mathbb{R}^3$

We can solve these problems using linear regression too

For multivariate problems, we use vector inputs

For multivariate problems, we use vector inputs

$$\boldsymbol{x} \in X; \quad X \in \mathbb{R}^3$$

For multivariate problems, we use vector inputs

$$\boldsymbol{x} \in X; \quad X \in \mathbb{R}^3$$

I will write

$$\boldsymbol{x}_i = \begin{bmatrix} x_i \langle 1 \rangle \\ x_i \langle 2 \rangle \\ x_i \langle 3 \rangle \end{bmatrix}$$

For multivariate problems, we use vector inputs

$$\boldsymbol{x} \in X; \quad X \in \mathbb{R}^3$$

I will write

$$\boldsymbol{x}_i = \begin{bmatrix} x_i\langle 1 \rangle \\ x_i\langle 2 \rangle \\ x_i\langle 3 \rangle \end{bmatrix}$$

$x_i\langle 1 \rangle$ refers to the first dimension of training data i

Assume an input space $X \in \mathbb{R}^\ell$

Assume an input space $X \in \mathbb{R}^\ell$

The design matrix for this **multivariate** linear system is

$$\mathbf{X}_D = [\mathbf{X} \ 1] = \begin{bmatrix} x_1 \langle \ell \rangle & x_1 \langle \ell - 1 \rangle & \dots & 1 \\ x_2 \langle \ell \rangle & x_2 \langle \ell - 1 \rangle & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ x_n \langle \ell \rangle & x_n \langle \ell - 1 \rangle & \dots & 1 \end{bmatrix}$$

Assume an input space $X \in \mathbb{R}^\ell$

The design matrix for this **multivariate** linear system is

$$X_D = [X \ 1] = \begin{bmatrix} x_1 \langle \ell \rangle & x_1 \langle \ell - 1 \rangle & \dots & 1 \\ x_2 \langle \ell \rangle & x_2 \langle \ell - 1 \rangle & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ x_n \langle \ell \rangle & x_n \langle \ell - 1 \rangle & \dots & 1 \end{bmatrix}$$

Remember $x_n \langle \ell \rangle$ refers to dimension ℓ of training data n

Assume an input space $X \in \mathbb{R}^\ell$

The design matrix for this **multivariate** linear system is

$$X_D = [X \ 1] = \begin{bmatrix} x_1 \langle \ell \rangle & x_1 \langle \ell - 1 \rangle & \dots & 1 \\ x_2 \langle \ell \rangle & x_2 \langle \ell - 1 \rangle & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ x_n \langle \ell \rangle & x_n \langle \ell - 1 \rangle & \dots & 1 \end{bmatrix}$$

Remember $x_n \langle \ell \rangle$ refers to dimension ℓ of training data n

We previously looked at linear and polynomial models for regression

We previously looked at linear and polynomial models for regression

$$f(\mathbf{x}, \boldsymbol{\theta}) = \mathbf{X}_D \boldsymbol{\theta} = \theta_m x^m + \theta_{m-1} x^{m-1} + \dots + \theta_0$$

We previously looked at linear and polynomial models for regression

$$f(\mathbf{x}, \boldsymbol{\theta}) = \mathbf{X}_D \boldsymbol{\theta} = \theta_m x^m + \theta_{m-1} x^{m-1} + \dots + \theta_0$$

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

Linear models are useful for certain problems

Linear models are useful for certain problems

1. Interpretability

Linear models are useful for certain problems

1. Interpretability
2. Low data requirement

Linear models are useful for certain problems

1. Interpretability
2. Low data requirement

But issues arise with other problems

Linear models are useful for certain problems

1. Interpretability
2. Low data requirement

But issues arise with other problems

1. Poor scalability

Linear models are useful for certain problems

1. Interpretability
2. Low data requirement

But issues arise with other problems

1. Poor scalability
2. Polynomials do not generalize well

Issues with very complex problems

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

Last time, we learned a polynomial function of a **one-dimensional x** using linear regression

Last time, we learned a polynomial function of a **one-dimensional** x using linear regression

However, we can also learn such functions for **multi-dimensional** x



$$X \in \mathbb{Z}_+^{256 \times 256}$$

Last time, we learned a polynomial function of a **one-dimensional** x using linear regression

However, we can also learn such functions for **multi-dimensional** x



$$X \in \mathbb{Z}_+^{256 \times 256}$$

This image contains 65536 pixels, so x has 65536 dimensions



256×256 pixels = 65536 pixels



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

What does the design matrix look like for an m-degree polynomial of this image?

$$X_D = \begin{bmatrix} x_1^m & x_1^{m-1} & \dots & x_1^1 & 1 \\ x_2^m & x_2^{m-1} & \dots & x_2^1 & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x_n^m & x_n^{m-1} & \dots & x_n^1 & 1 \\ x_1^{n-1}x_2 & x^{n-2}x_2^2 & \dots & 0 & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x_1x_2\dots x_n & 0 & \dots & 0 & 1 \end{bmatrix}$$

$$X_D = \begin{bmatrix} x_1^m & x_1^{m-1} & \dots & x_1^1 & 1 \\ x_2^m & x_2^{m-1} & \dots & x_2^1 & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x_n^m & x_n^{m-1} & \dots & x_n^1 & 1 \\ x_1^{n-1}x_2 & x^{n-2}x_2^2 & \dots & 0 & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x_1x_2\dots x_n & 0 & \dots & 0 & 1 \end{bmatrix}$$

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

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

$$X_D = \begin{bmatrix} x_1^m & x_1^{m-1} & \dots & x_1^1 & 1 \\ x_2^m & x_2^{m-1} & \dots & x_2^1 & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x_n^m & x_n^{m-1} & \dots & x_n^1 & 1 \\ x_1^{n-1}x_2 & x^{n-2}x_2^2 & \dots & 0 & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x_1x_2\dots x_n & 0 & \dots & 0 & 1 \end{bmatrix}$$

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

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

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

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

We must invert $\mathbf{X}_D^\top \mathbf{X}_D$, 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

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

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

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

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

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

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

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

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

One day this will be possible, but today it is not

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

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

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

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

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

One day this will be possible, but today it is not

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**

What happens to polynomials outside of the support (dataset)?

What happens to polynomials outside of the support (dataset)?

$$\theta_m x^m + \theta_{m-1} x^{m-1} + \dots$$

Equation of a polynomial

What happens to polynomials outside of the support (dataset)?

$$\theta_m x^m + \theta_{m-1} x^{m-1} + \dots$$

Equation of a polynomial

$$x^m \left(\theta_m + \frac{\theta_{m-1}}{x} + \dots \right)$$

Factor out x^m

What happens to polynomials outside of the support (dataset)?

$$\theta_m x^m + \theta_{m-1} x^{m-1} + \dots$$

Equation of a polynomial

$$x^m \left(\theta_m + \frac{\theta_{m-1}}{x} + \dots \right)$$

Factor out x^m

$$\lim_{x \rightarrow \infty} x^m \left(\theta_m + \frac{\theta_{m-1}}{x} + \dots \right)$$

Take the limit

What happens to polynomials outside of the support (dataset)?

$$\theta_m x^m + \theta_{m-1} x^{m-1} + \dots$$

Equation of a polynomial

$$x^m \left(\theta_m + \frac{\theta_{m-1}}{x} + \dots \right)$$

Factor out x^m

$$\lim_{x \rightarrow \infty} x^m \left(\theta_m + \frac{\theta_{m-1}}{x} + \dots \right)$$

Take the limit

$$\lim_{x \rightarrow \infty} x^m \cdot \lim_{x \rightarrow \infty} \left(\theta_m + \frac{\theta_{m-1}}{x} + \dots \right)$$

Split the limit (limit of products)

$$\lim_{x \rightarrow \infty} x^m \cdot \lim_{x \rightarrow \infty} \left(\theta_m + \frac{\theta_{m-1}}{x} + \dots \right) \quad \text{Split the limit (limit of products)}$$

$$\left(\lim_{x \rightarrow \infty} x^m \right) \cdot (\theta_m + 0 + \dots) \quad \text{Evaluate right limit}$$

$$\lim_{x \rightarrow \infty} x^m \cdot \lim_{x \rightarrow \infty} \left(\theta_m + \frac{\theta_{m-1}}{x} + \dots \right) \quad \text{Split the limit (limit of products)}$$

$$\left(\lim_{x \rightarrow \infty} x^m \right) \cdot (\theta_m + 0 + \dots) \quad \text{Evaluate right limit}$$

$$\theta_m \lim_{x \rightarrow \infty} x^m \quad \text{Rewrite}$$

$$\lim_{x \rightarrow \infty} x^m \cdot \lim_{x \rightarrow \infty} \left(\theta_m + \frac{\theta_{m-1}}{x} + \dots \right) \quad \text{Split the limit (limit of products)}$$

$$\left(\lim_{x \rightarrow \infty} x^m \right) \cdot (\theta_m + 0 + \dots) \quad \text{Evaluate right limit}$$

$$\theta_m \lim_{x \rightarrow \infty} x^m \quad \text{Rewrite}$$

$$\theta_m \lim_{x \rightarrow \infty} x^m = \infty \quad \text{If } \theta_m > 0$$

$$\lim_{x \rightarrow \infty} x^m \cdot \lim_{x \rightarrow \infty} \left(\theta_m + \frac{\theta_{m-1}}{x} + \dots \right) \quad \text{Split the limit (limit of products)}$$

$$\left(\lim_{x \rightarrow \infty} x^m \right) \cdot (\theta_m + 0 + \dots) \quad \text{Evaluate right limit}$$

$$\theta_m \lim_{x \rightarrow \infty} x^m \quad \text{Rewrite}$$

$$\theta_m \lim_{x \rightarrow \infty} x^m = \infty \quad \text{If } \theta_m > 0$$

$$\theta_m \lim_{x \rightarrow \infty} x^m = -\infty \quad \text{If } \theta_m < 0$$

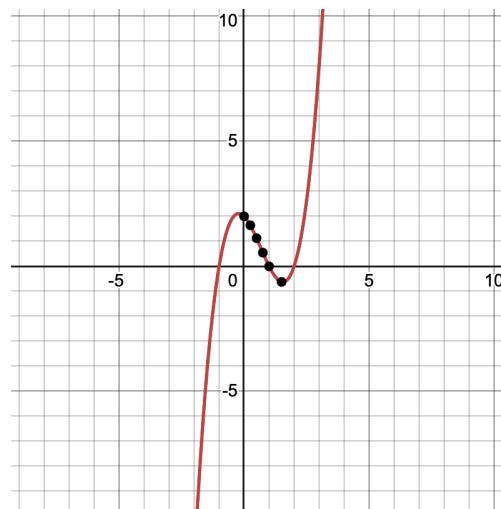
Polynomials quickly tend towards $-\infty, \infty$ outside of the support

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

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

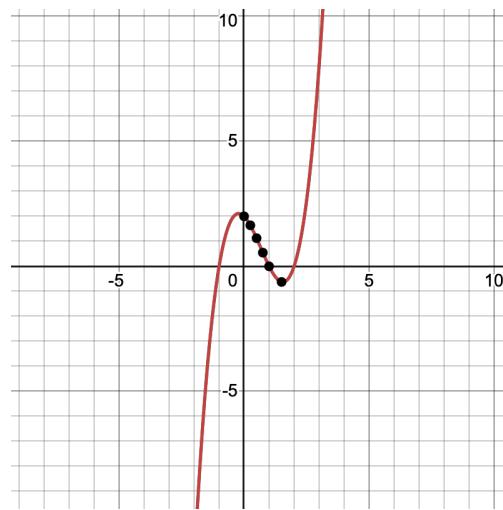
Polynomials quickly tend towards $-\infty, \infty$ outside of the support

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



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

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



Remember, to predict new data we want our functions to generalize

Linear regression has issues

Linear regression has issues

1. Poor scalability

Linear regression has issues

1. Poor scalability
2. Polynomials do not generalize well

1. Review
2. **Multivariate linear regression**
3. Limitations of linear regression
4. History of neural networks
5. Biological neurons
6. Artificial neurons
7. Perceptron
8. Multilayer Perceptron

1. Review
2. Multivariate linear regression
3. **Limitations of linear regression**
4. History of neural networks
5. Biological neurons
6. Artificial neurons
7. Perceptron
8. Multilayer Perceptron

Can we improve upon linear regression?

Can we improve upon linear regression?

Yes, with neural networks

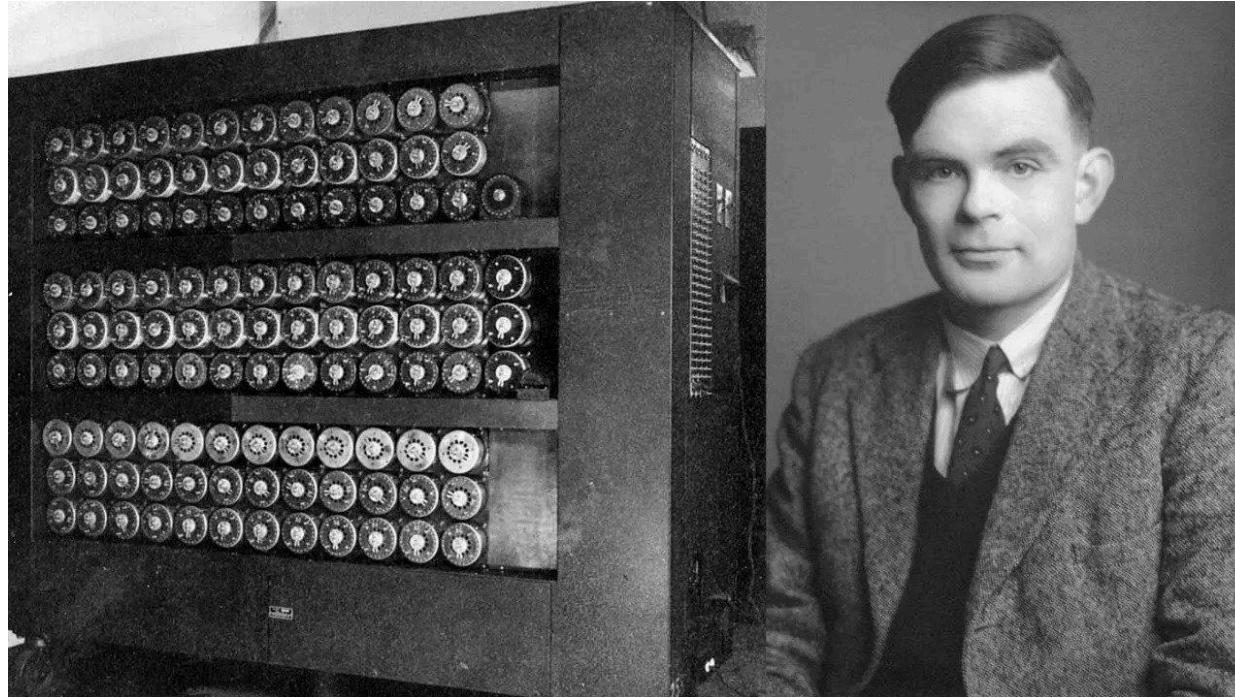
In 1939-1945, there was a World War

In 1939-1945, there was a World War

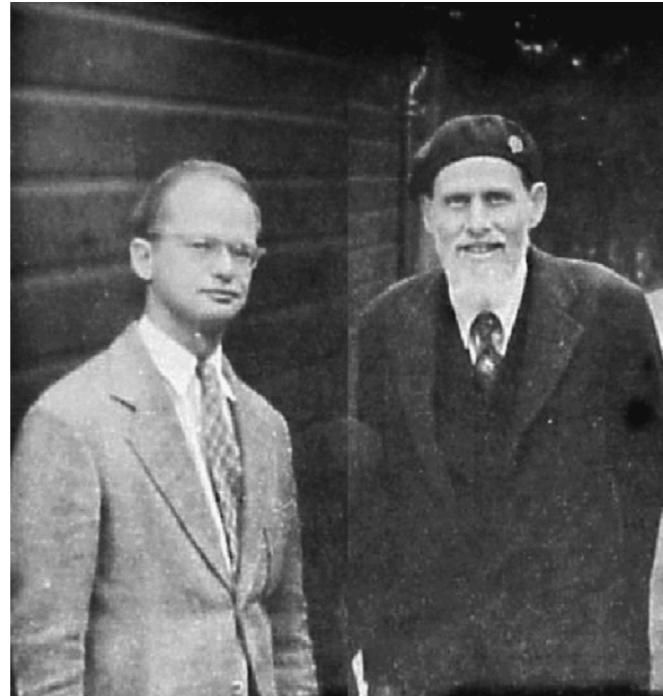
Militaries invested funding for research, and invented the computer

In 1939-1945, there was a World War

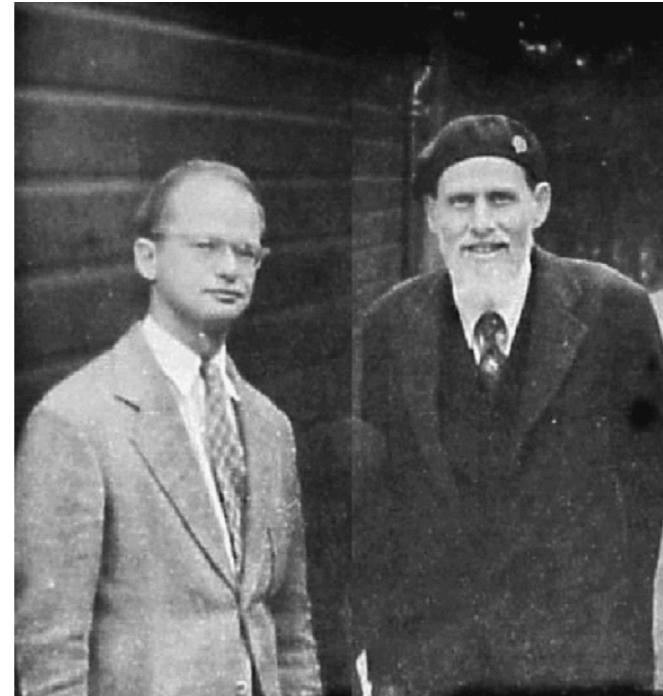
Militaries invested funding for research, and invented the computer



Meanwhile, a neuroscientist and mathematician (McCullough and Pitts) were trying to understand the human brain

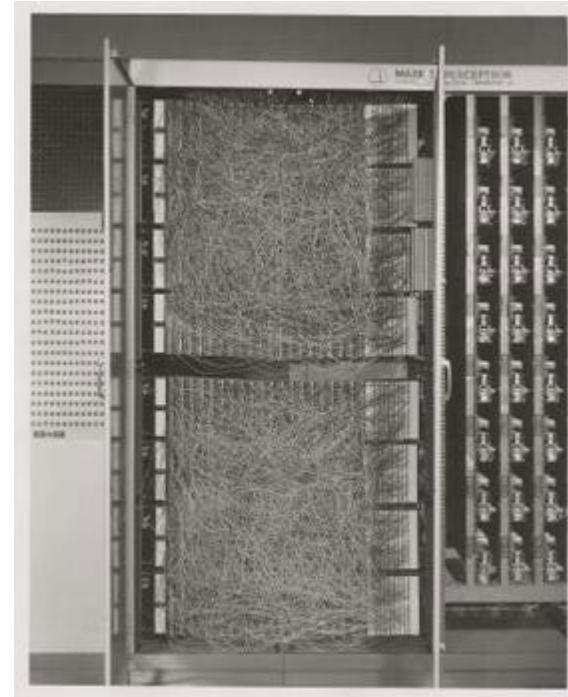


Meanwhile, a neuroscientist and mathematician (McCullough and Pitts) were trying to understand the human brain



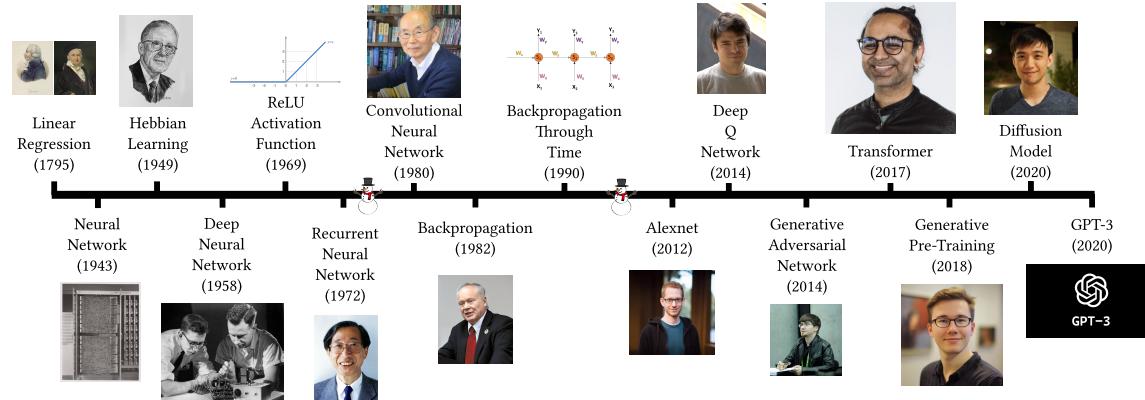
They designed the theory for the first neural network

A few years later, Rosenblatt implemented this neural network using a new invention – the computer

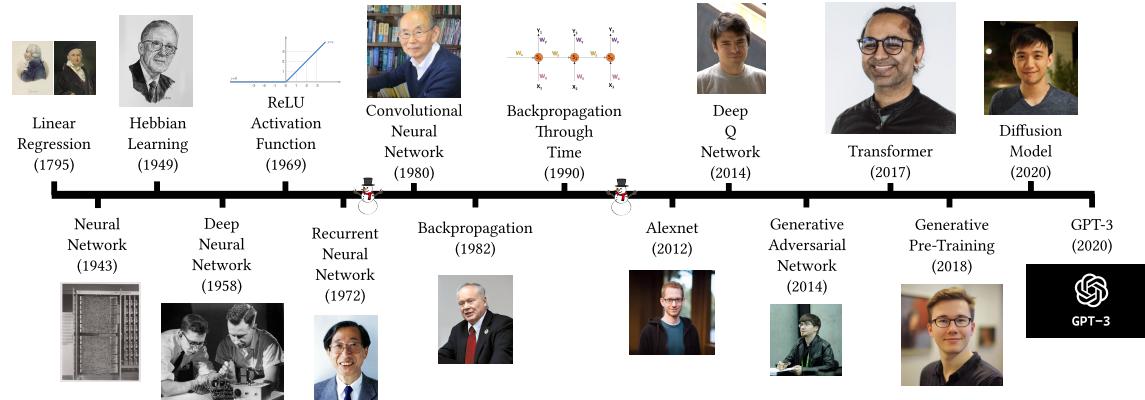


Through advances in theory and hardware, neural networks became slightly better

Through advances in theory and hardware, neural networks became slightly better



Through advances in theory and hardware, neural networks became slightly better



Around 2012, these improvements culminated in neural networks that perform like humans

So what is a neural network?

So what is a neural network?

It is a function, inspired by how the brain works

So what is a neural network?

It is a function, inspired by how the brain works

$$f : X \times \Theta \mapsto Y$$

Brains and neural networks rely on **neurons**

Brains and neural networks rely on **neurons**

Brain: Biological neurons → Biological neural network

Brains and neural networks rely on **neurons**

Brain: Biological neurons → Biological neural network

Computer: Artificial neurons → Artificial neural network

Brains and neural networks rely on **neurons**

Brain: Biological neurons → Biological neural network

Computer: Artificial neurons → Artificial neural network

First, let us review biological neurons

Brains and neural networks rely on **neurons**

Brain: Biological neurons → Biological neural network

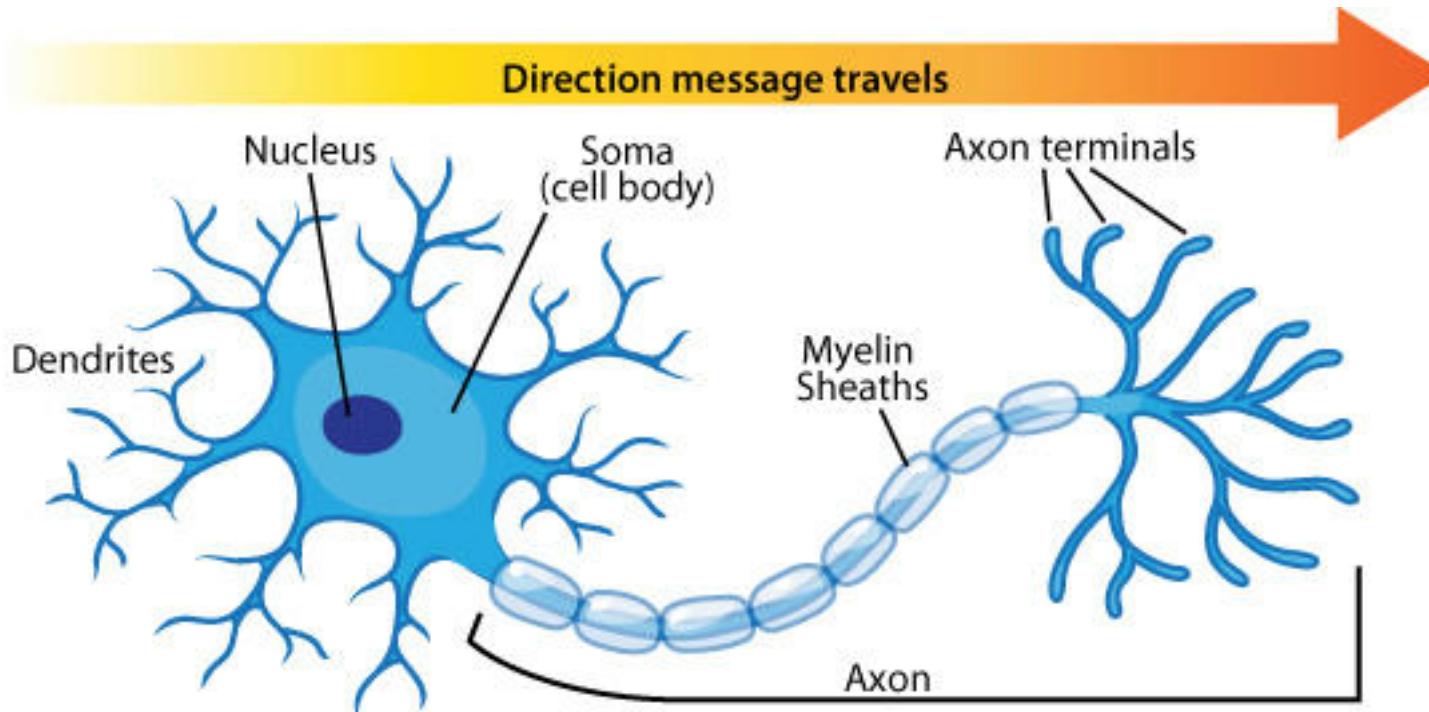
Computer: Artificial neurons → Artificial neural network

First, let us review biological neurons

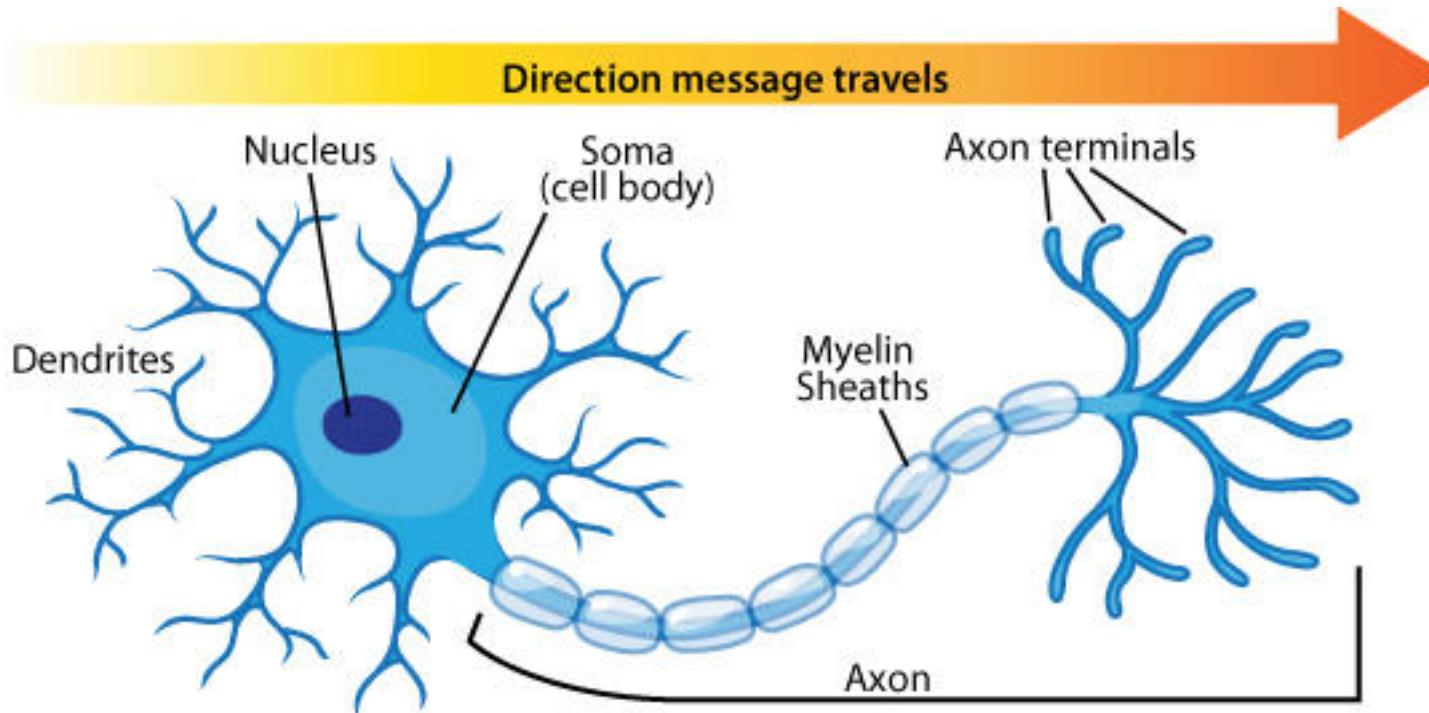
Note: I am not a neuroscientist! I may make simplifications or errors with biology

1. Review
2. Multivariate linear regression
3. **Limitations of linear regression**
4. History of neural networks
5. Biological neurons
6. Artificial neurons
7. Perceptron
8. Multilayer Perceptron

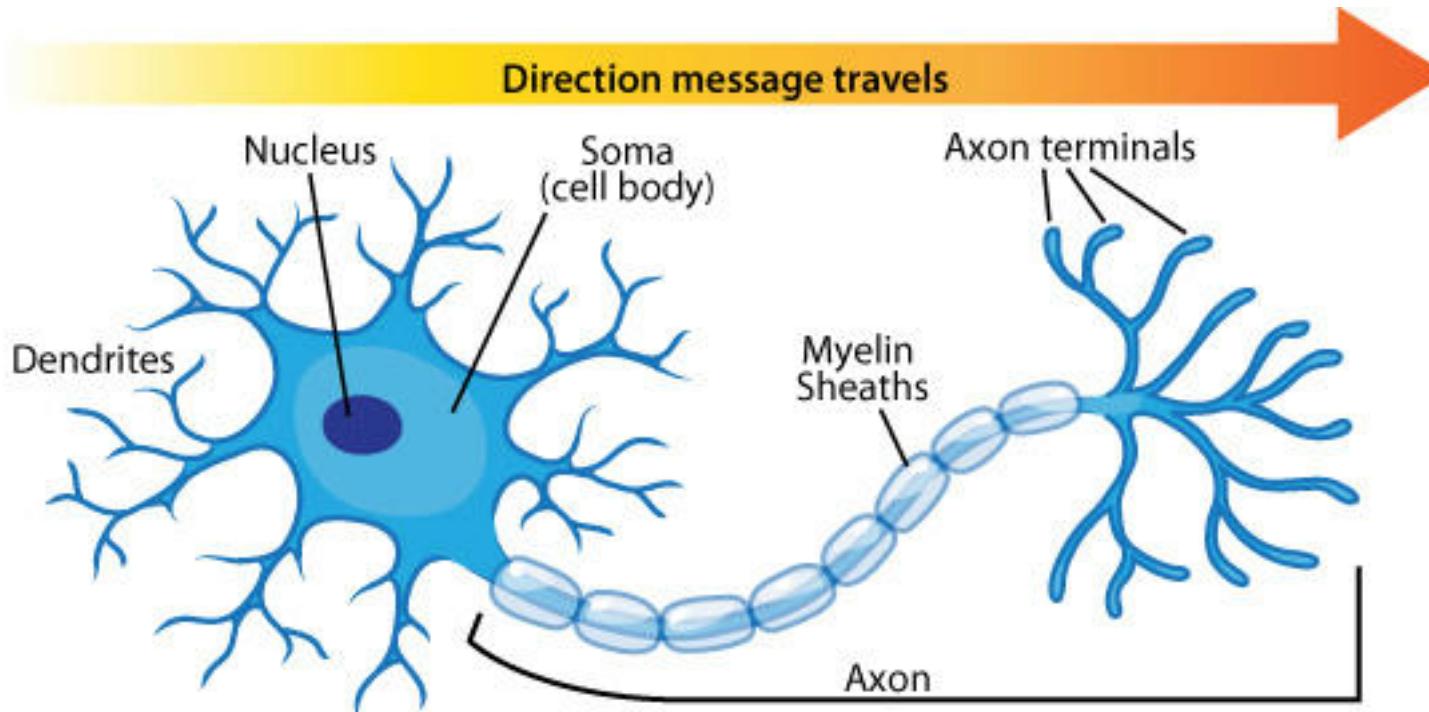
1. Review
2. Multivariate linear regression
3. Limitations of linear regression
4. **History of neural networks**
5. Biological neurons
6. Artificial neurons
7. Perceptron
8. Multilayer Perceptron



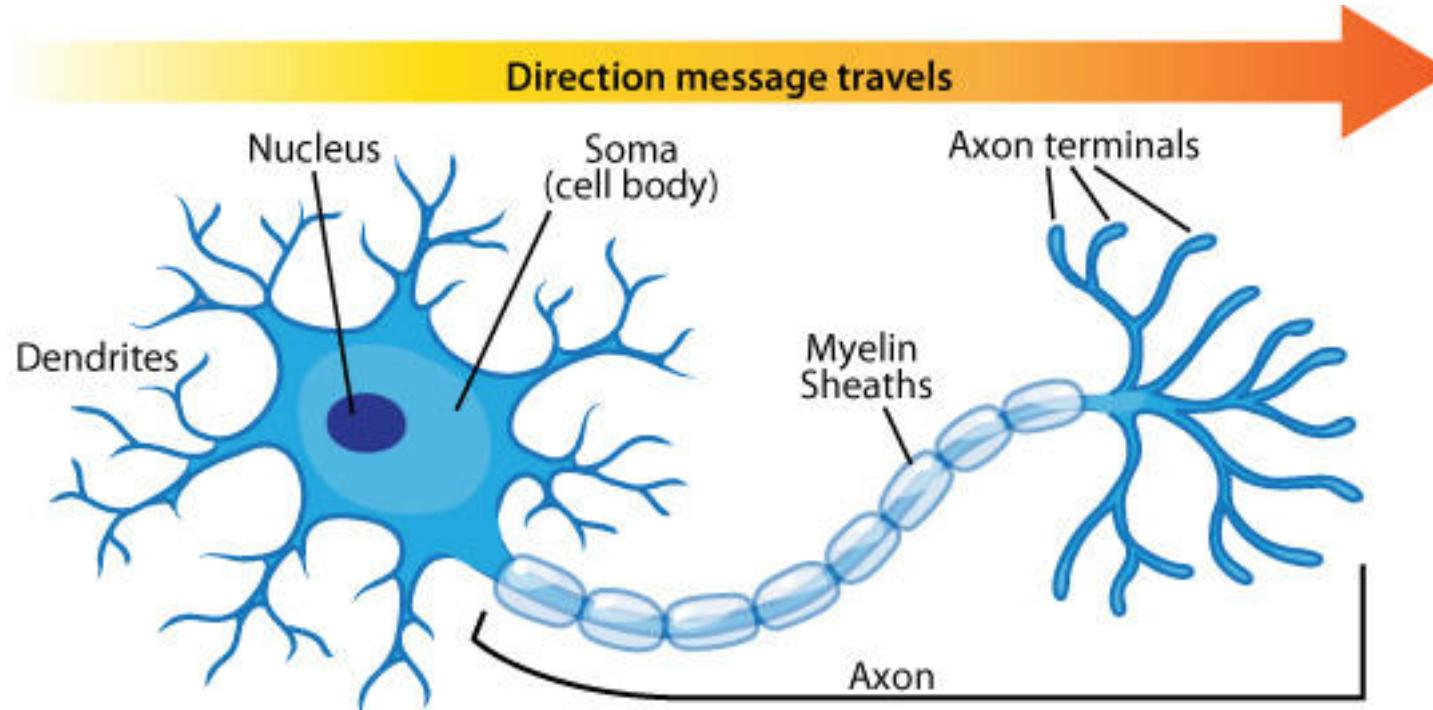
A simplified neuron consists of many parts



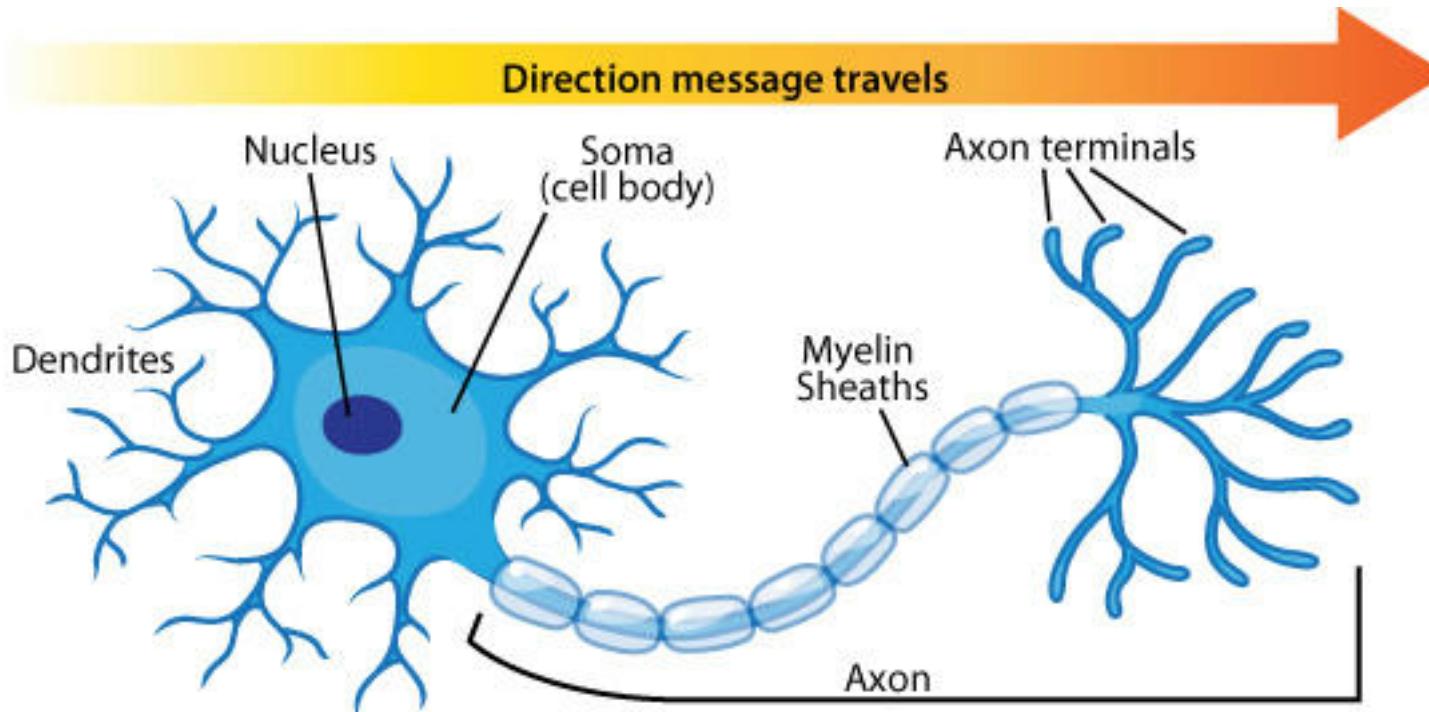
Neurons send and process messages from other neurons



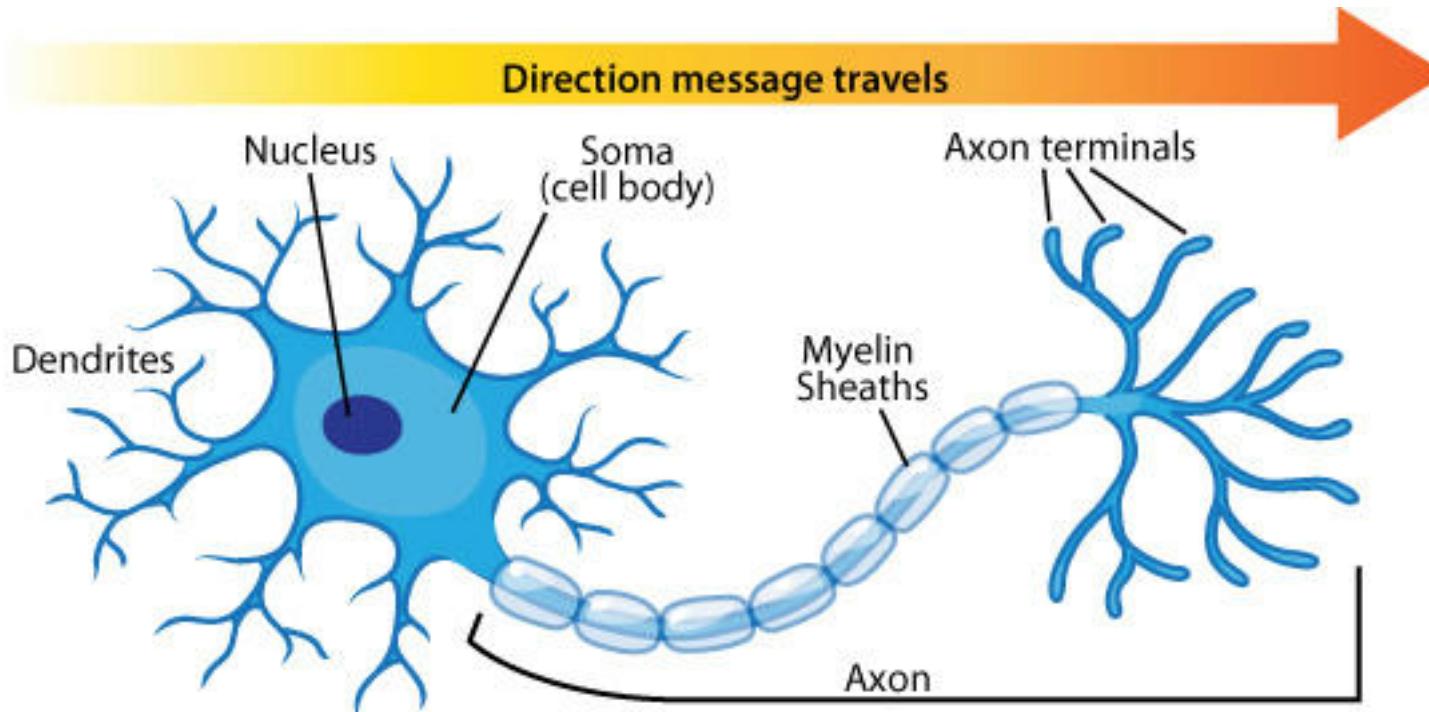
Incoming electrical signals travel along dendrites



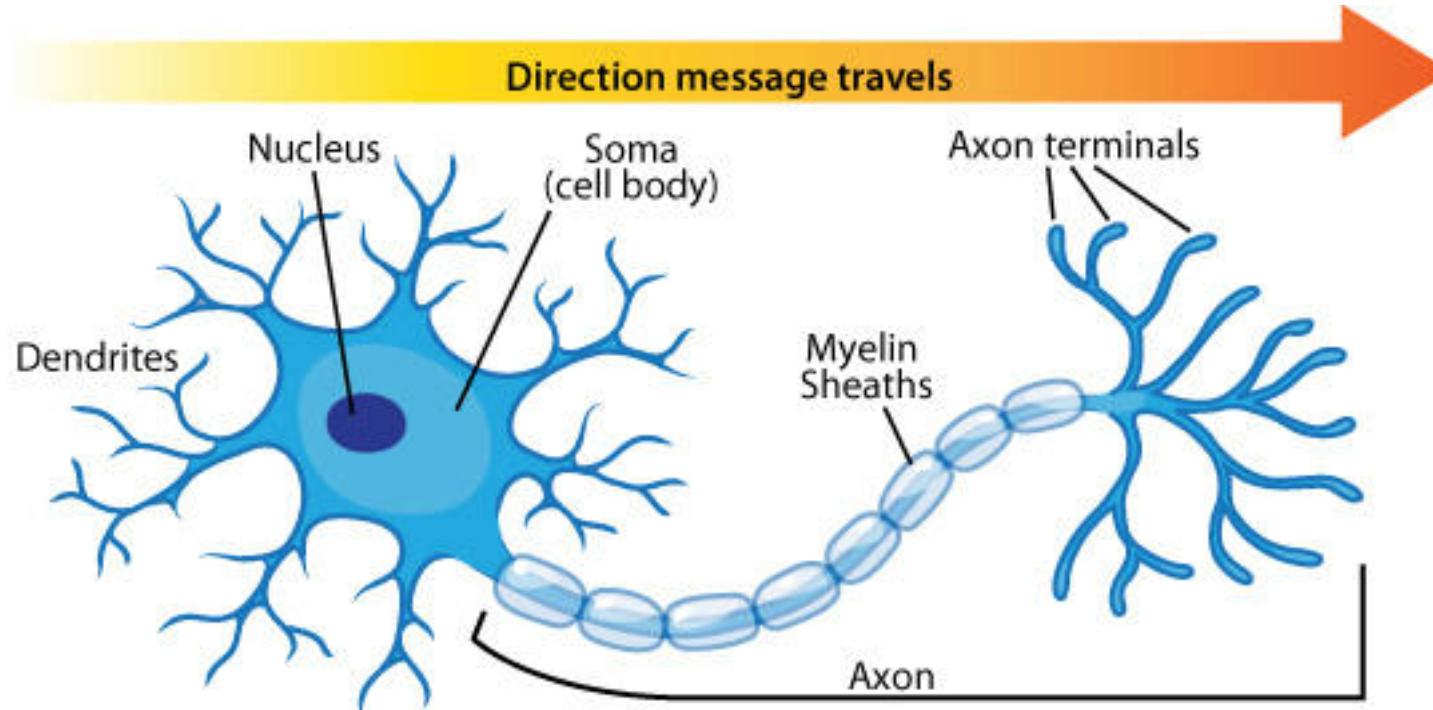
Dendrites are not all equal! Different dendrites have different diameters and structures



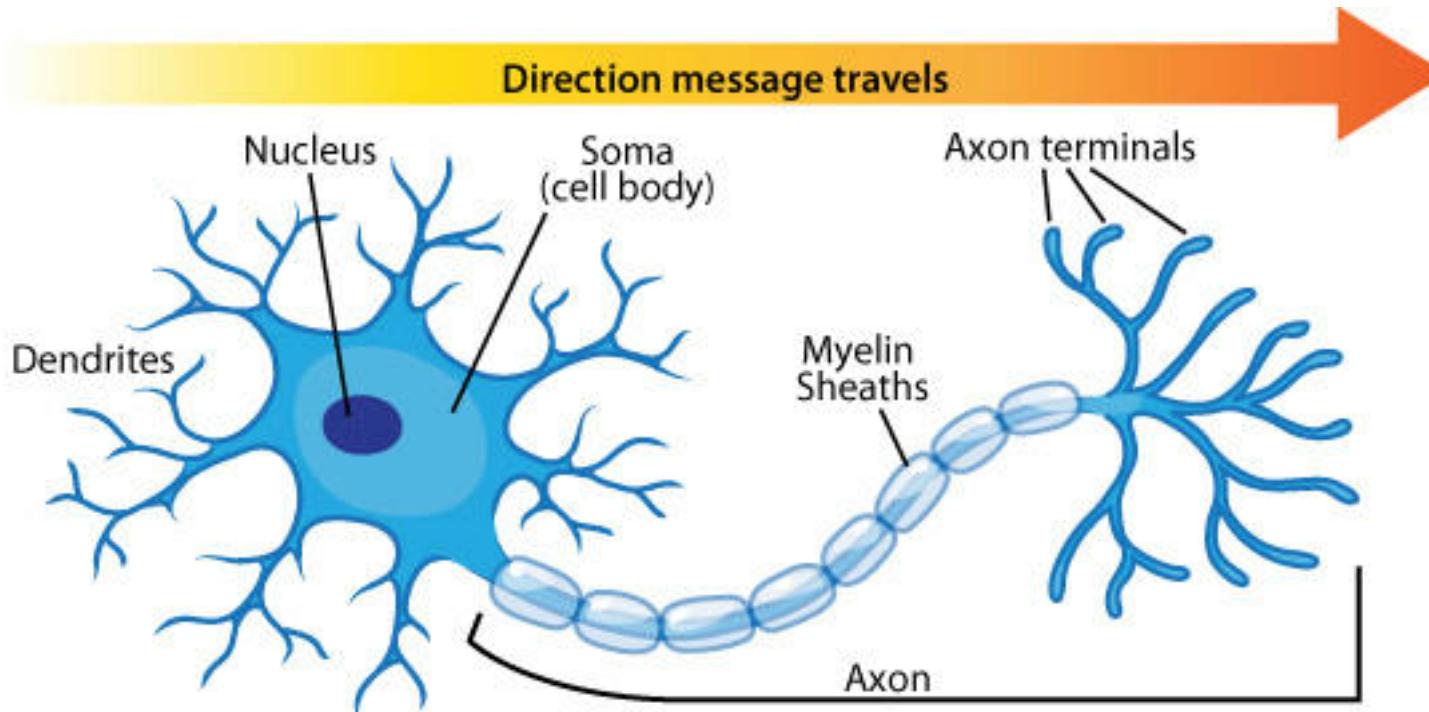
The axon outputs an electrical signal to other neurons



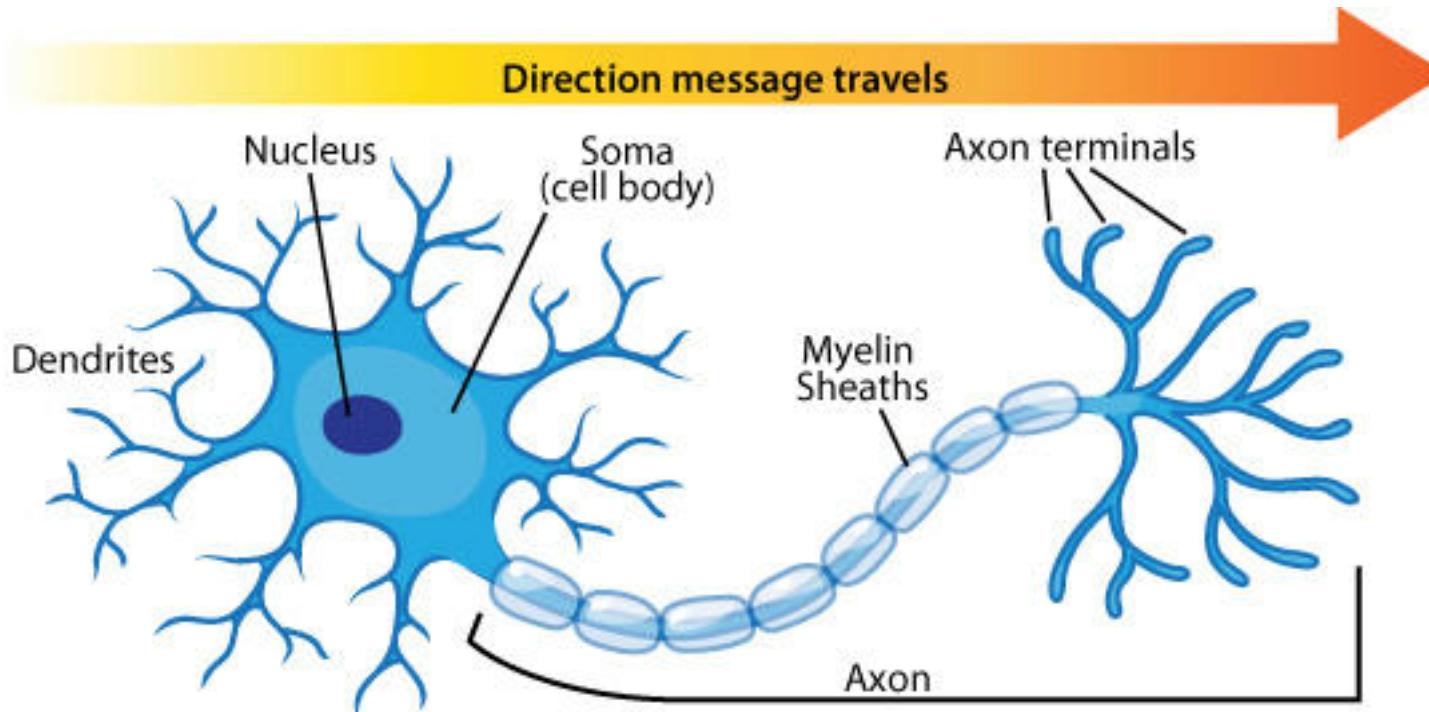
The axon terminals will connect to dendrites of other neurons



For our purposes, we can consider the axon terminals and dendrites to be the same thing



The neuron takes many inputs, and produces a single output

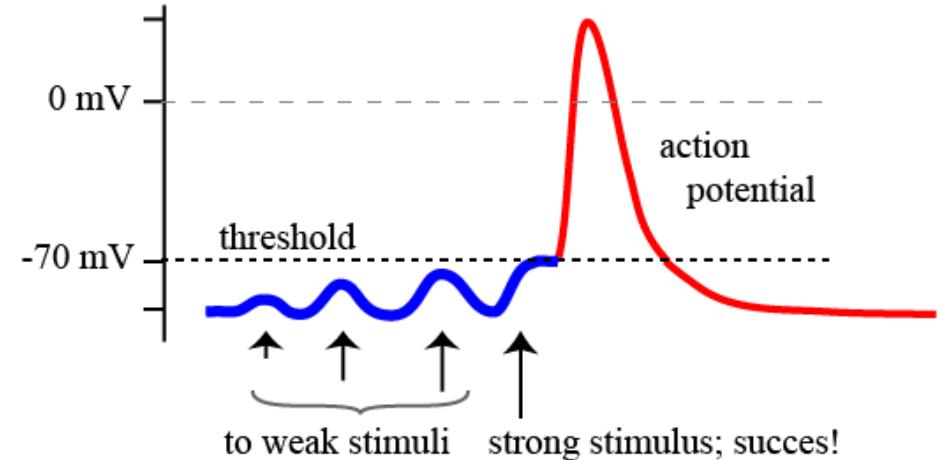


The neuron will only output a signal down the axon (“fire”) at certain times

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

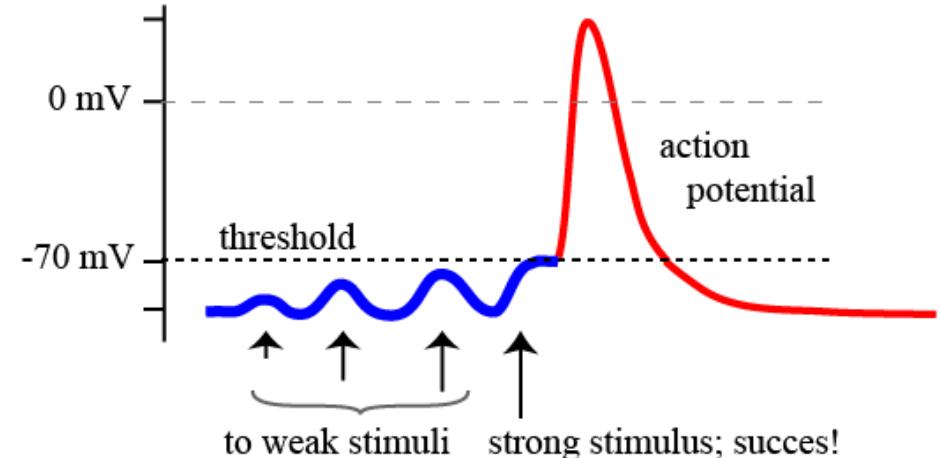
How does a neuron decide to send an impulse (“fire”)?

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



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

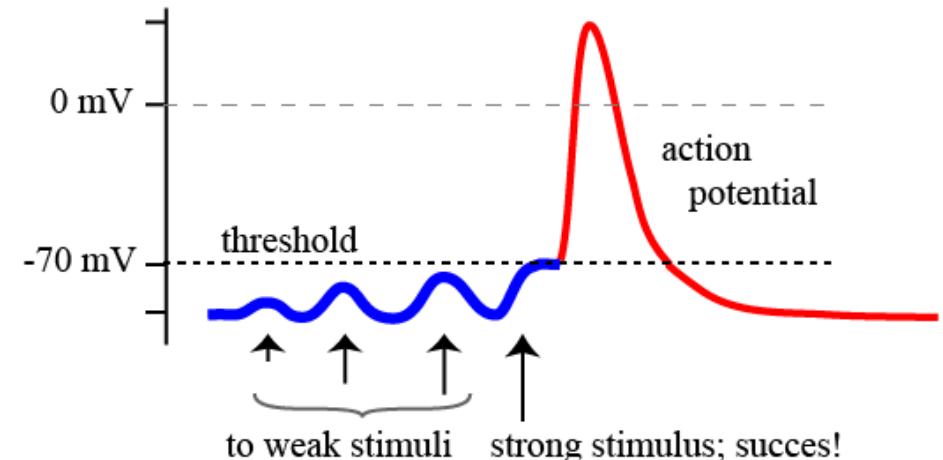
Incoming impulses (via dendrites)
change the electric potential of the
neuron



Recall that in a parallel circuit, we can sum voltages together

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

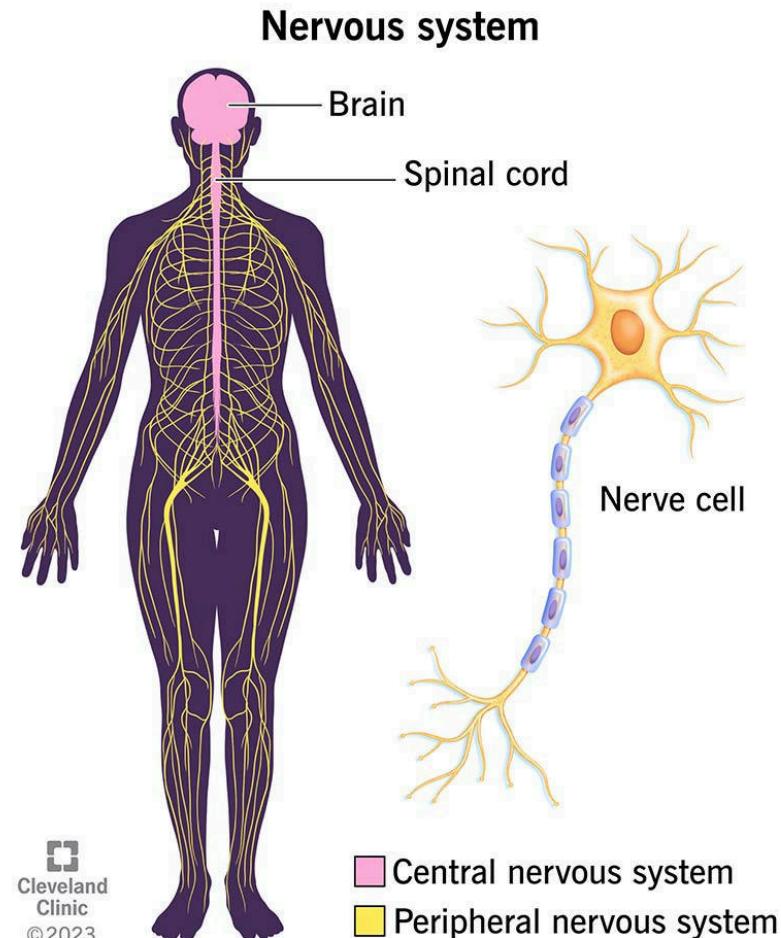
Incoming impulses (via dendrites) change the electric potential of the neuron



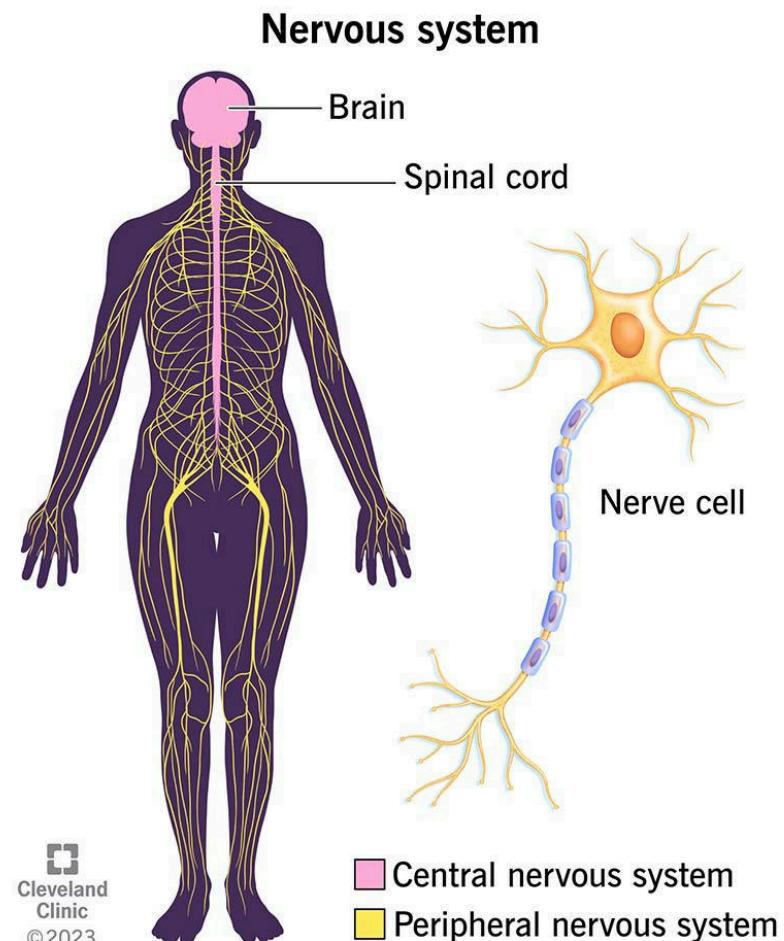
Recall that in a parallel circuit, we can sum voltages together

Many active dendrites will add together and trigger an impulse

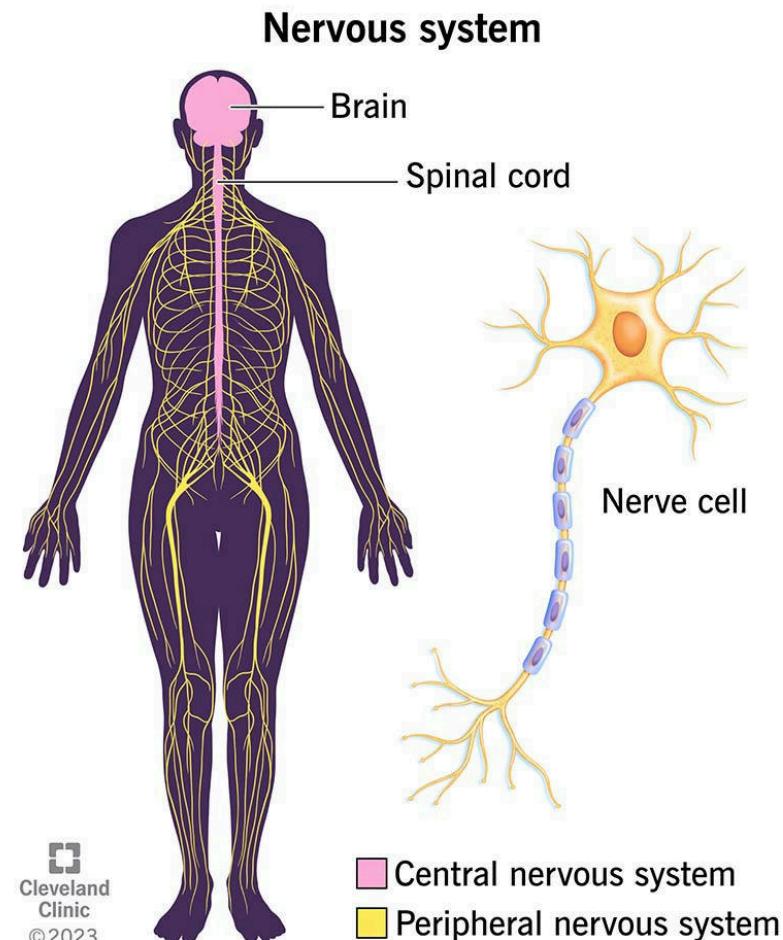
Pain triggers initial nerve impulse,
starts a chain reaction into the
brain



When the signal reaches the brain,
we will think

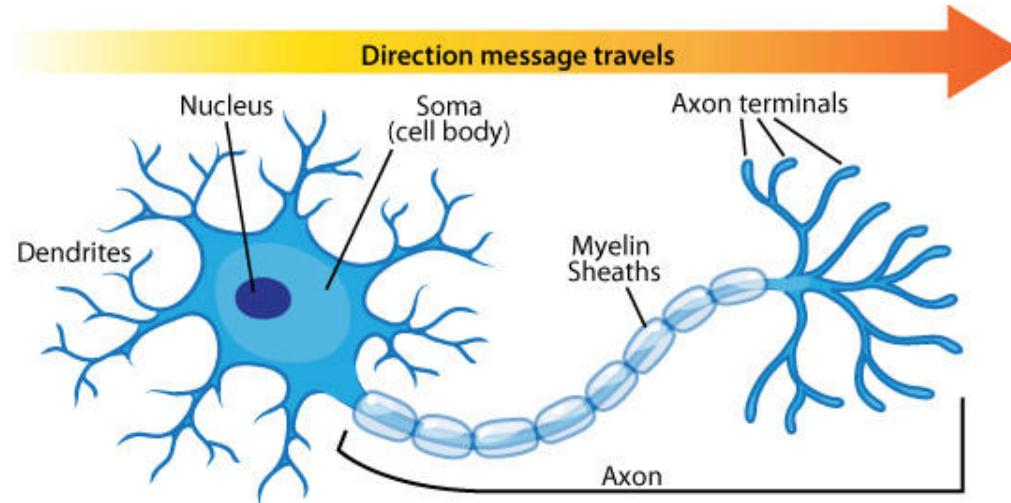


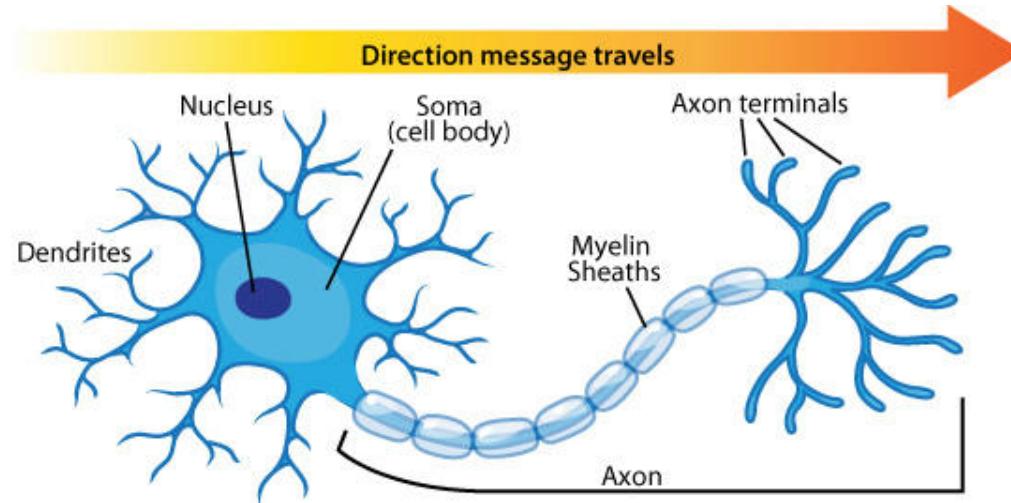
After thinking, we will take action



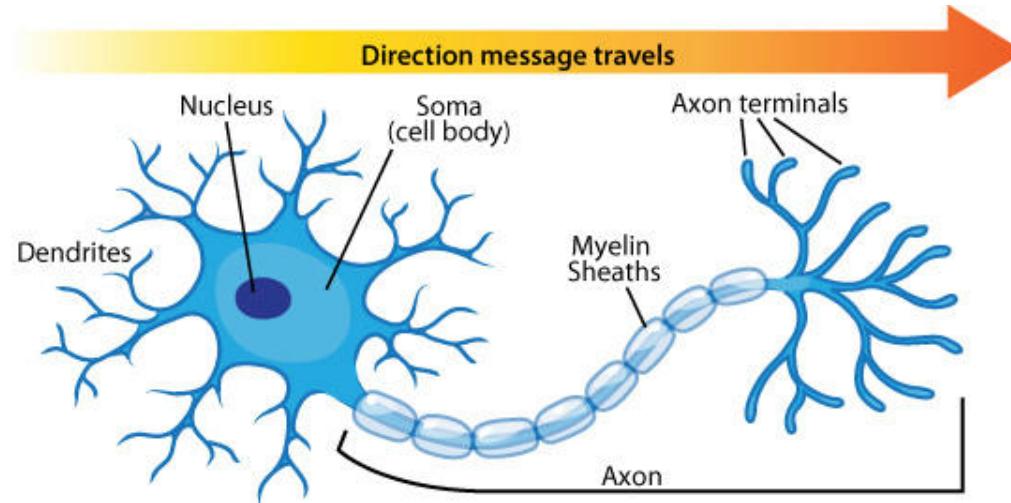
1. Review
2. Multivariate linear regression
3. Limitations of linear regression
4. **History of neural networks**
5. Biological neurons
6. Artificial neurons
7. Perceptron
8. Multilayer Perceptron

1. Review
2. Multivariate linear regression
3. Limitations of linear regression
4. History of neural networks
5. **Biological neurons**
6. Artificial neurons
7. Perceptron
8. Multilayer Perceptron





Question: How could we write a neuron as a function? $f : \underline{\quad} \mapsto \underline{\quad}$



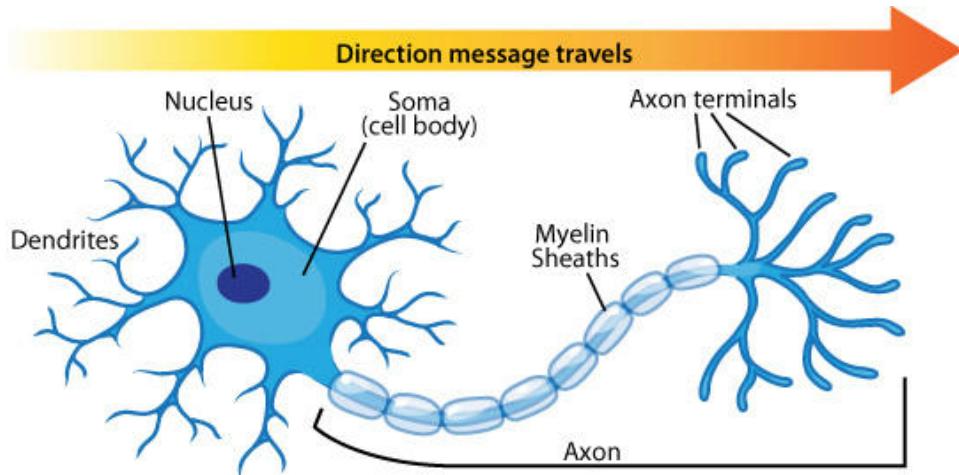
Question: How could we write a neuron as a function? $f : \underline{\quad} \mapsto \underline{\quad}$

Answer:

$$f : \underbrace{\mathbb{R}^m}_{\text{Dendrite voltages}} \times \underbrace{\mathbb{R}^m}_{\text{Dendrite size}} \mapsto \underbrace{\mathbb{R}}_{\text{Axon voltage}}$$

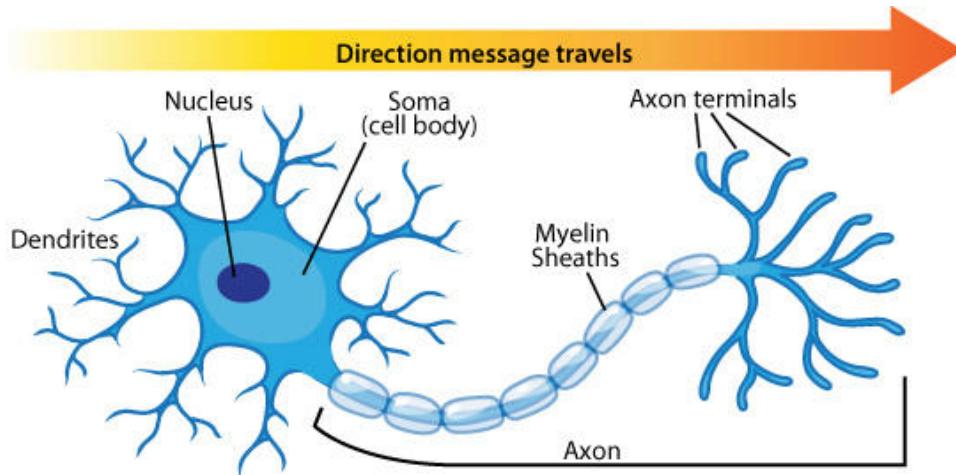
Let us implement an artifical neuron as a function

Let us implement an artifical neuron as a function



Neuron has a structure of
dendrites

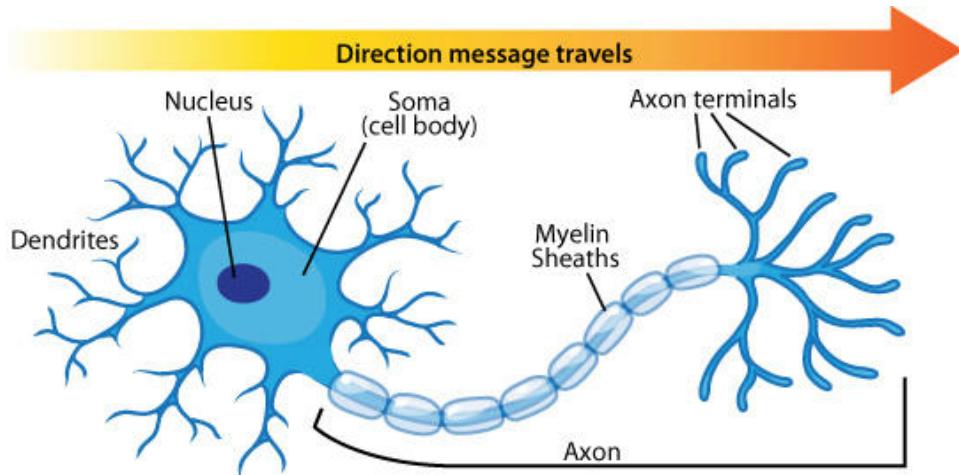
Let us implement an artifical neuron as a function



Neuron has a structure of dendrites

$$f \left(\begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \right) = f \left(\begin{bmatrix} 0.5 \\ 3.1 \\ \vdots \\ 2.0 \end{bmatrix} \right)$$

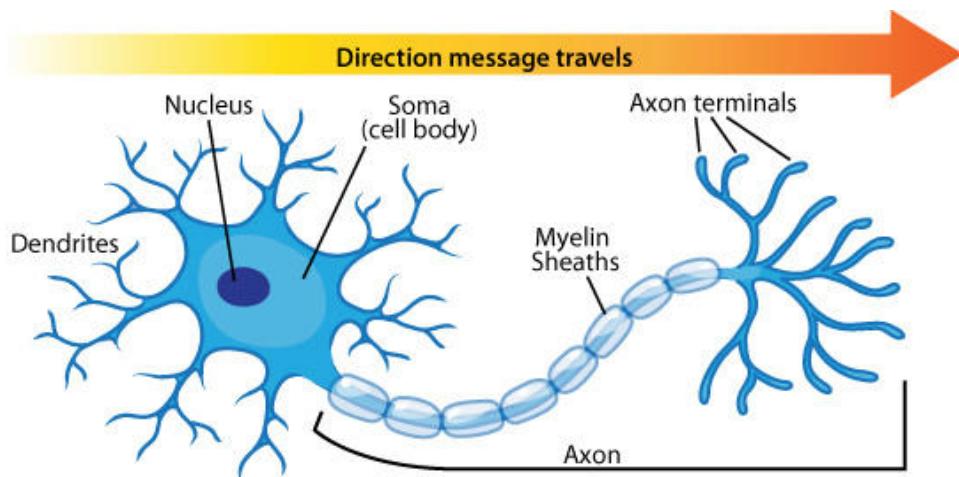
Let us implement an artifical neuron as a function



Each incoming dendrite has some voltage potential

Let us implement an artificial neuron as a function

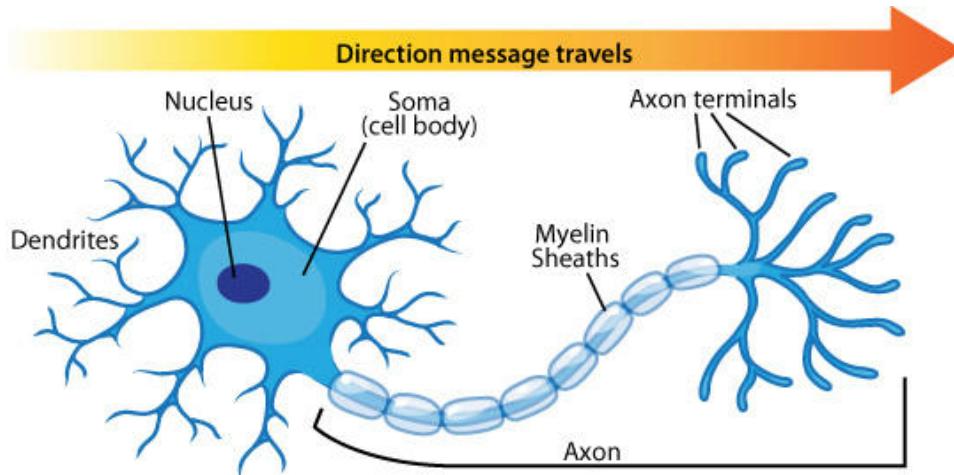
Each incoming dendrite has some voltage potential



$$f \left(\begin{bmatrix} x_i\langle 1 \rangle \\ \vdots \\ x_i\langle n \rangle \end{bmatrix}, \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \right)$$

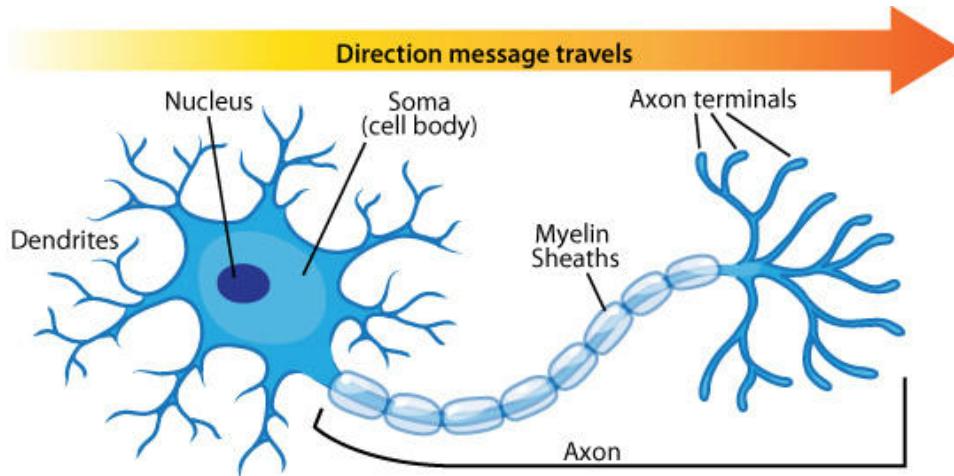
$$x_i = \begin{bmatrix} x_i\langle 1 \rangle \\ \vdots \\ x_i\langle n \rangle \end{bmatrix} = \begin{bmatrix} 0.5 \\ \vdots \\ -0.3 \end{bmatrix}$$

Let us implement an artificial neuron as a function



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

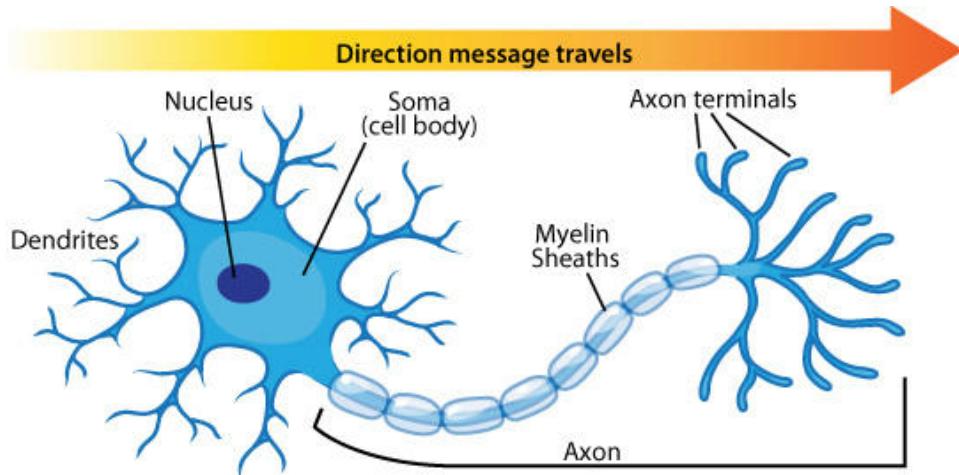
Let us implement an artificial neuron as a function



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

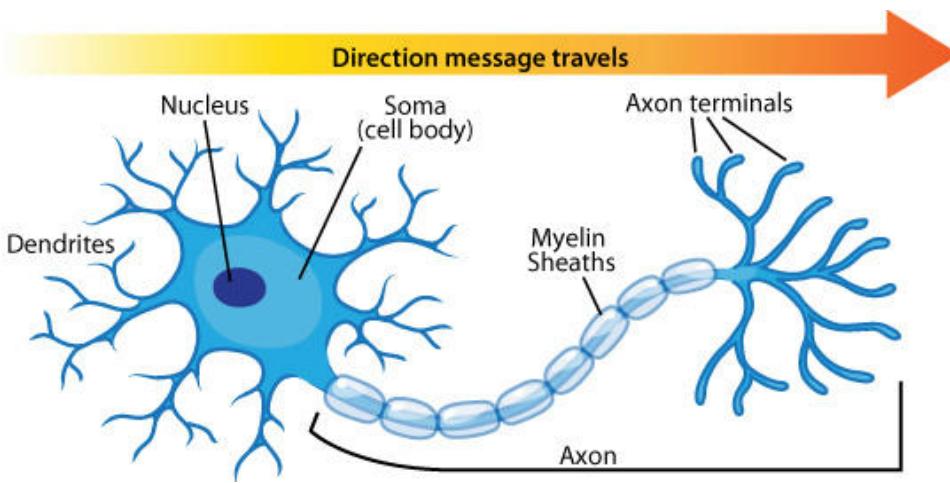
$$f \left(\begin{bmatrix} x_i \langle 1 \rangle \\ \vdots \\ x_i \langle n \rangle \end{bmatrix}, \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \right) = \sum_{j=1}^n x_i \langle j \rangle \theta_j$$

Let us implement an artificial neuron as a function



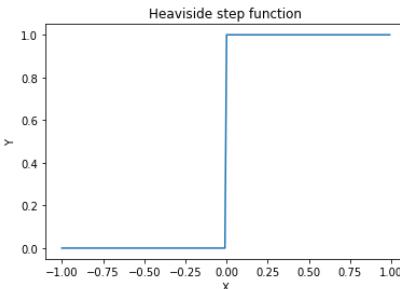
The axon fires only if the voltage
is over a threshold

Let us implement an artificial neuron as a function



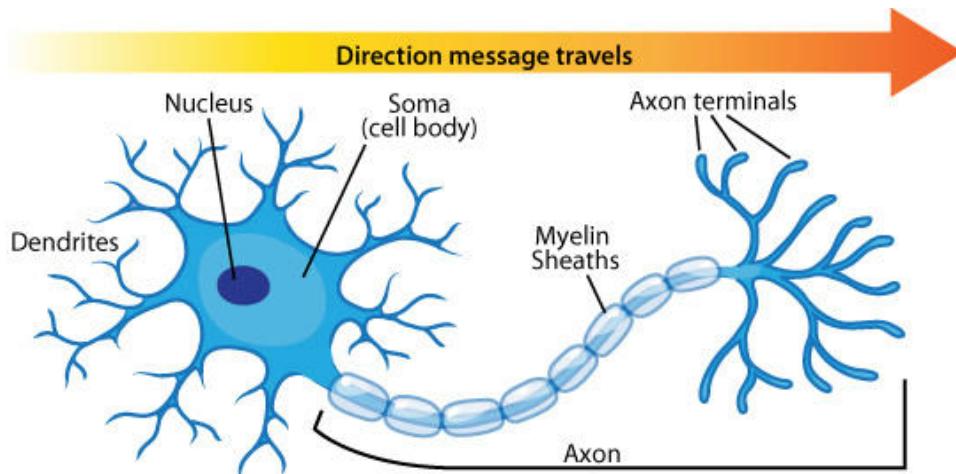
The axon fires only if the voltage
is over a threshold

$$\sigma(x) =$$

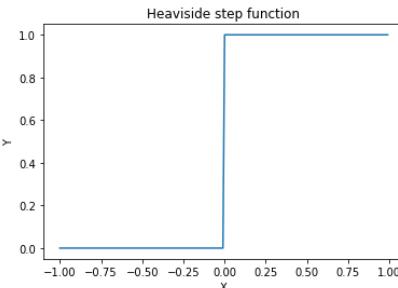


Let us implement an artifical neuron as a function

The axon fires only if the voltage
is over a threshold



$$\sigma(x) =$$



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

This is almost the artificial neuron!

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

This is almost the artificial neuron!

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

$$f(\boldsymbol{x}, \boldsymbol{\theta}) = \sigma\left(\sum_{j=1}^n x \langle j \rangle \theta_j\right)$$

This is almost the artificial neuron!

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

$$f(\boldsymbol{x}, \boldsymbol{\theta}) = \sigma\left(\sum_{j=1}^n x \langle j \rangle \theta_j\right)$$

Let us write this out for clarity

$$f(\mathbf{x}, \boldsymbol{\theta}) = \sigma(x\langle n \rangle \theta_n + x\langle n - 1 \rangle \theta_{n-1} + \dots + x\langle 1 \rangle \theta_1)$$

$$f(\mathbf{x}, \boldsymbol{\theta}) = \sigma(x\langle n \rangle \theta_n + x\langle n - 1 \rangle \theta_{n-1} + \dots + x\langle 1 \rangle \theta_1)$$

Question: Does this look familiar to anyone?

$$f(\mathbf{x}, \boldsymbol{\theta}) = \sigma(x\langle n \rangle \theta_n + x\langle n - 1 \rangle \theta_{n-1} + \dots + x\langle 1 \rangle \theta_1)$$

Question: Does this look familiar to anyone?

Answer: This is a multivariate linear model!

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

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

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

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

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

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

Artificial neuron

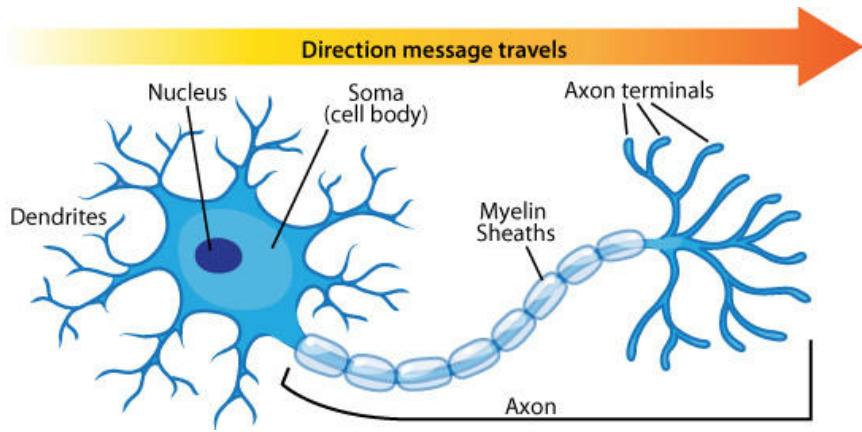
$$f(\mathbf{x}, \boldsymbol{\theta}) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

Linear model

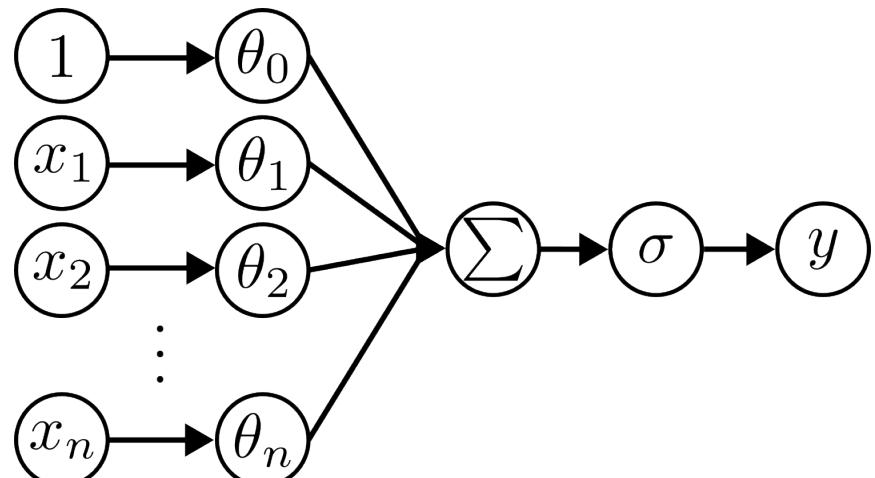
It is the linear model with an activation function!

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

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



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



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

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

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

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

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

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

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

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

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

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

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

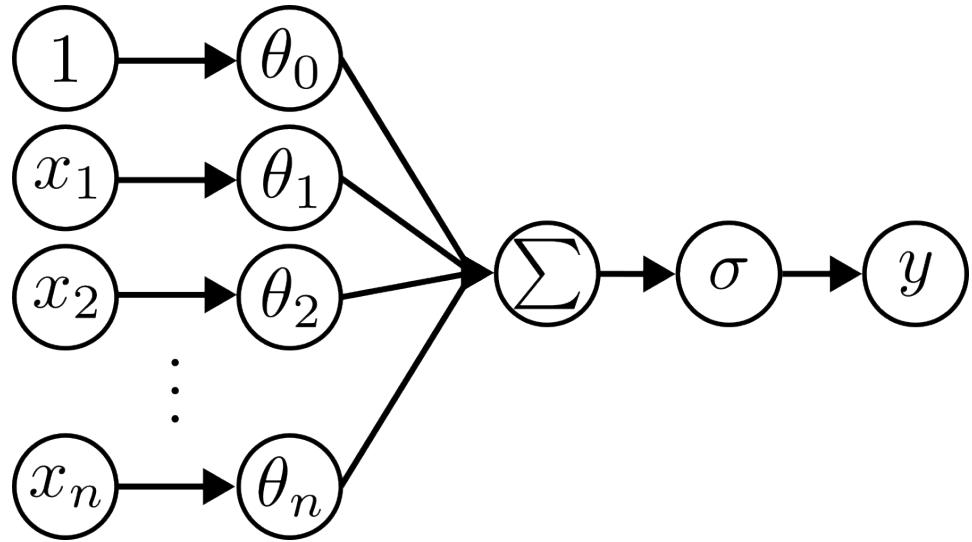
$$f\left(\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}, \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_n \end{bmatrix}\right) = \sigma(\theta_0 + \theta_{1:n} \cdot \mathbf{x})$$

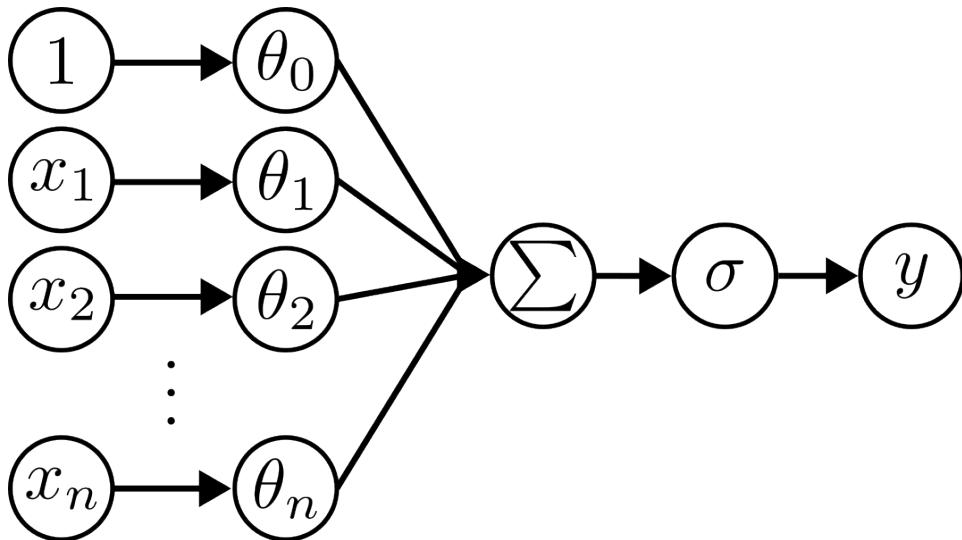
We often write the parameters as a **weight** w and **bias** b

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

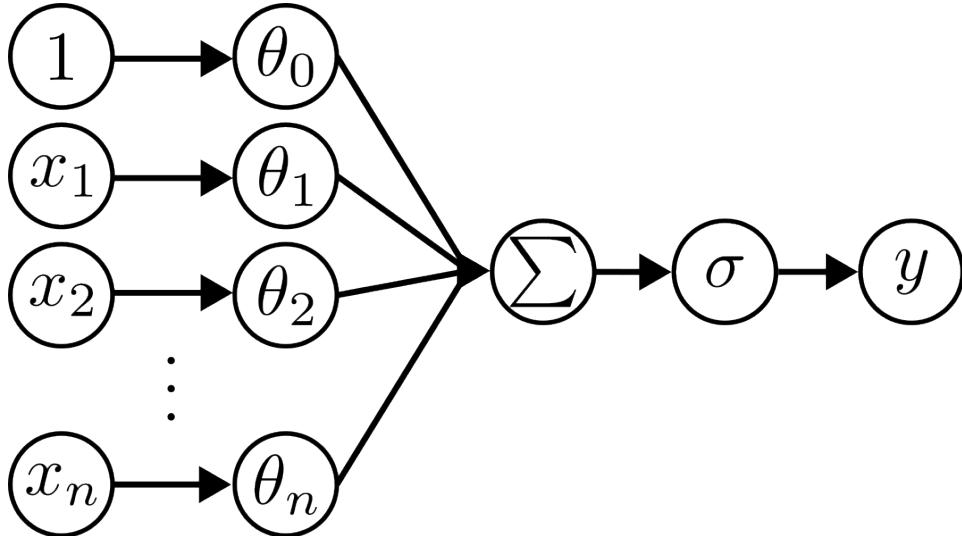
$$b = \theta_0, \mathbf{w} = \theta_{1:n}$$

Relax



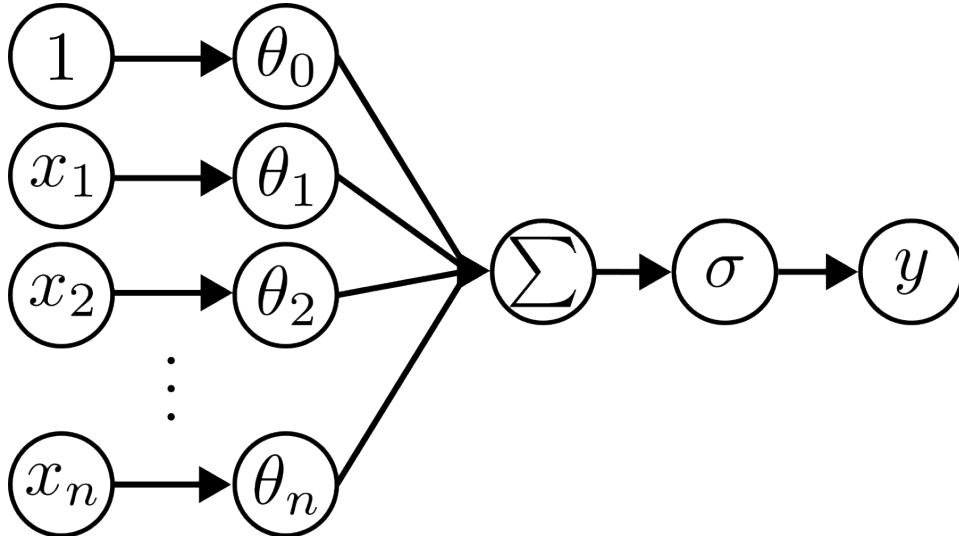


Recall that in machine learning we deal with functions



Recall that in machine learning we deal with functions

What kinds of functions can our neuron represent?

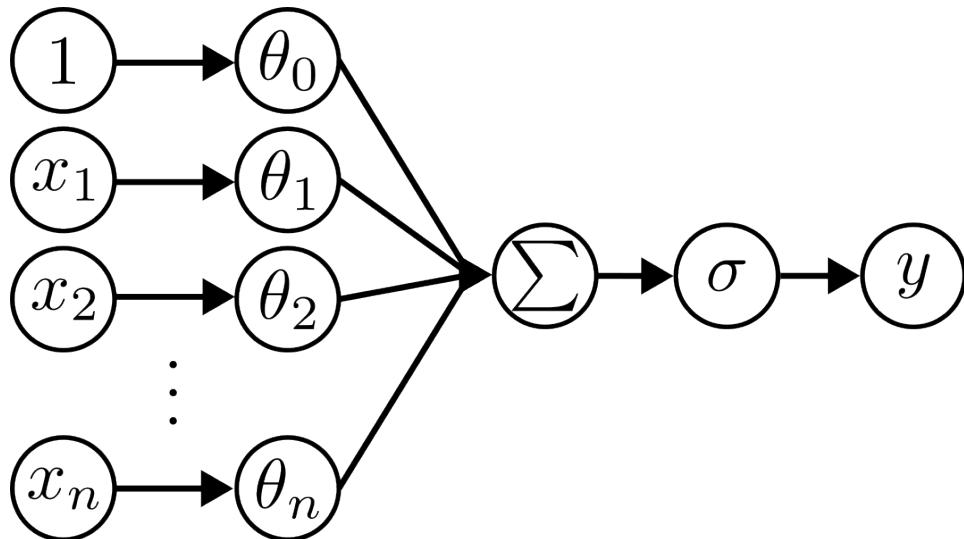


Recall that in machine learning we deal with functions

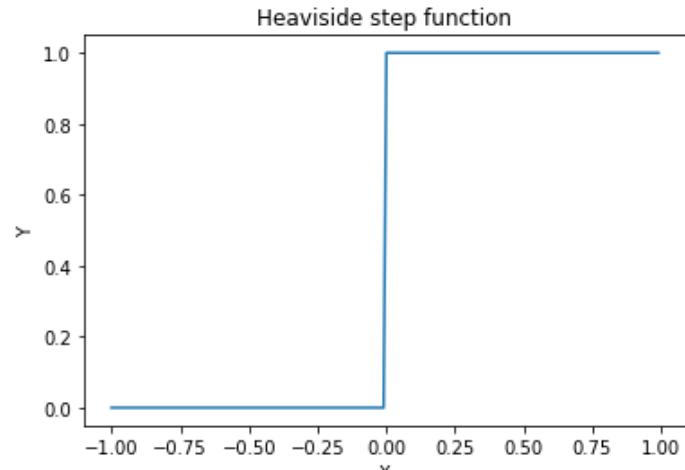
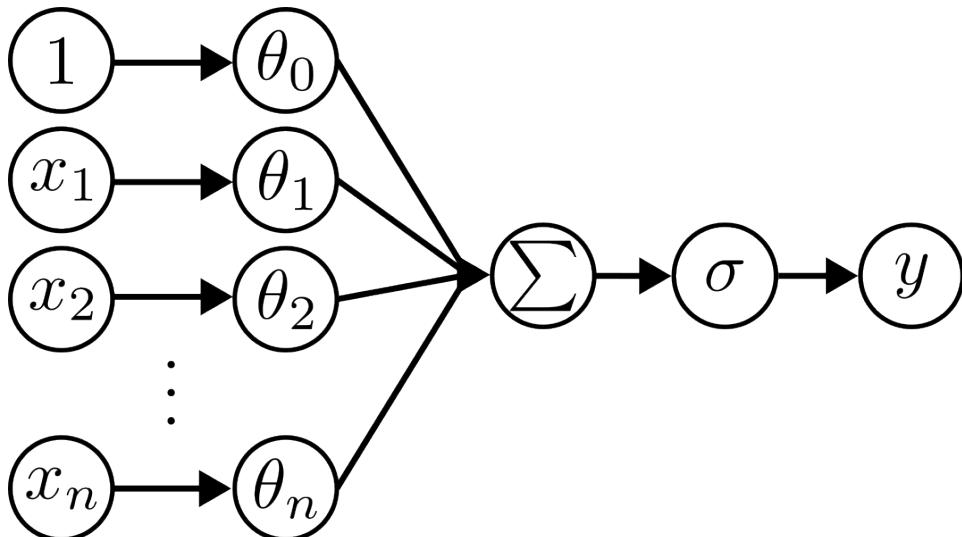
What kinds of functions can our neuron represent?

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

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

Implement AND using an artificial neuron

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

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

Implement AND using an artificial neuron

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

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

x_1	x_2	y	$f(x_1, x_2, \theta)$	\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

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

Implement OR using an artificial neuron

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

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

Implement OR using an artificial neuron

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

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

x_1	x_2	y	$f(x_1, x_2, \theta)$	\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

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

Implement XOR using an artificial neuron

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

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

Implement XOR using an artificial neuron

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

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

x_1	x_2	y	$f(x_1, x_2, \theta)$	\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?

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

Why can't we represent XOR using a neuron?

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

We can only represent H (linear function)

Why can't we represent XOR using a neuron?

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

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?

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

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?

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

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

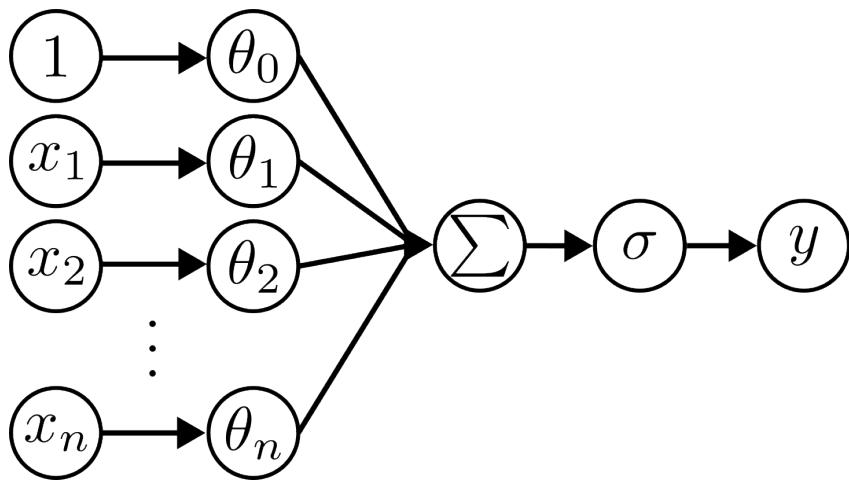
Let us think back to biology, maybe it has an answer

Brain: Biological neurons → Biological neural network

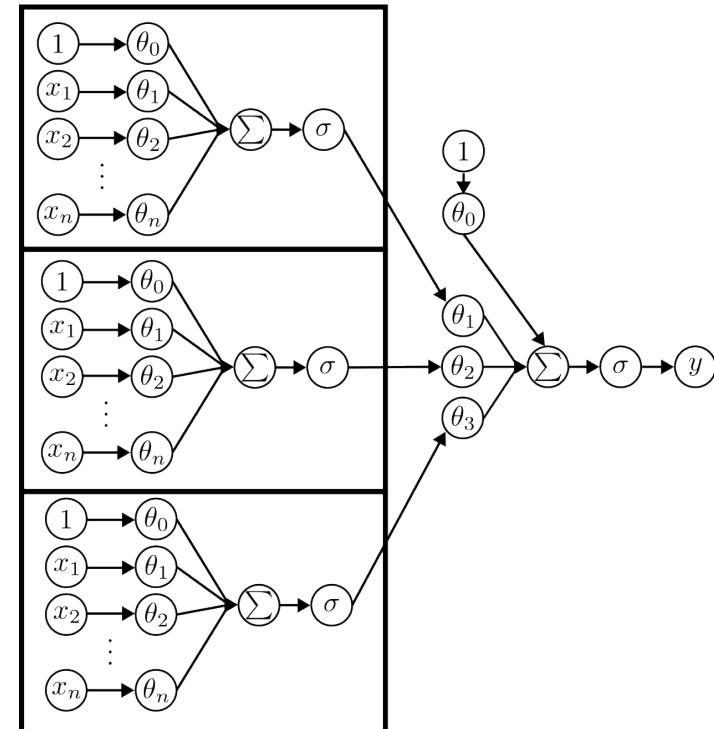
Brain: Biological neurons → Biological neural network

Computer: Artificial neurons → Artificial neural network

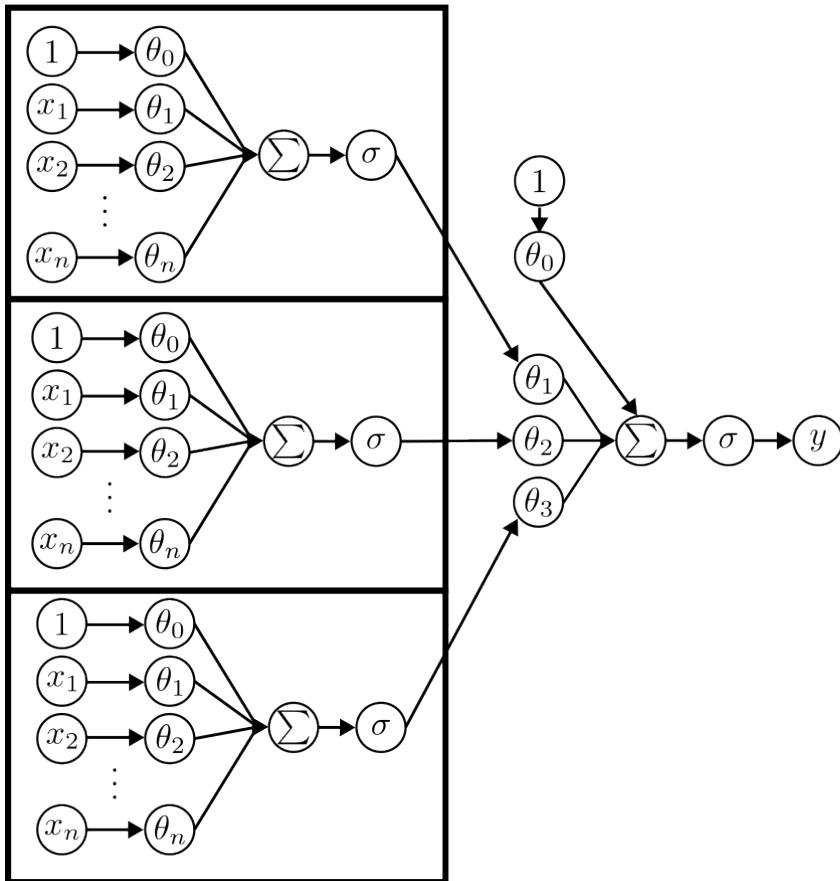
Connect artificial neurons into a network



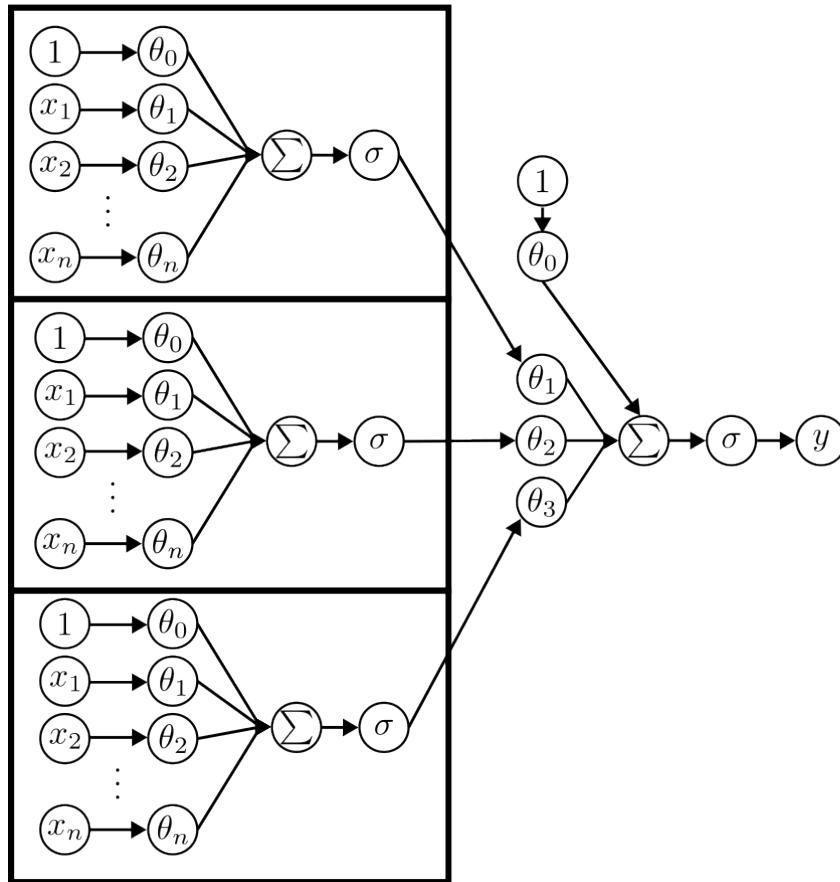
Neuron



Neural Network

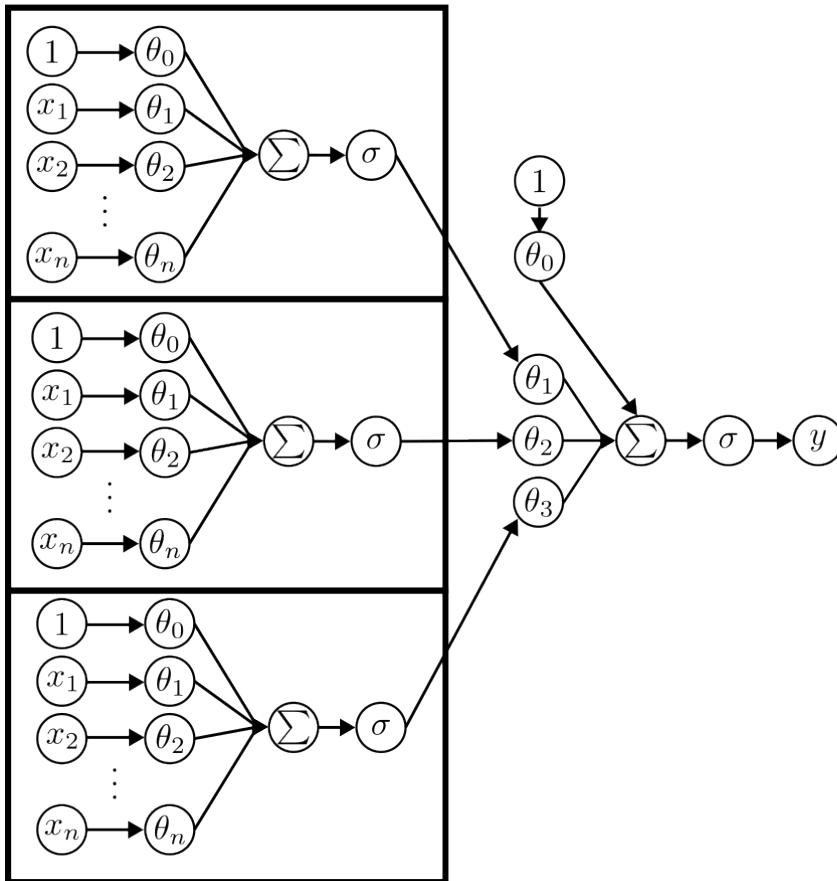


Adding neurons in **parallel**
creates a **wide** neural network



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

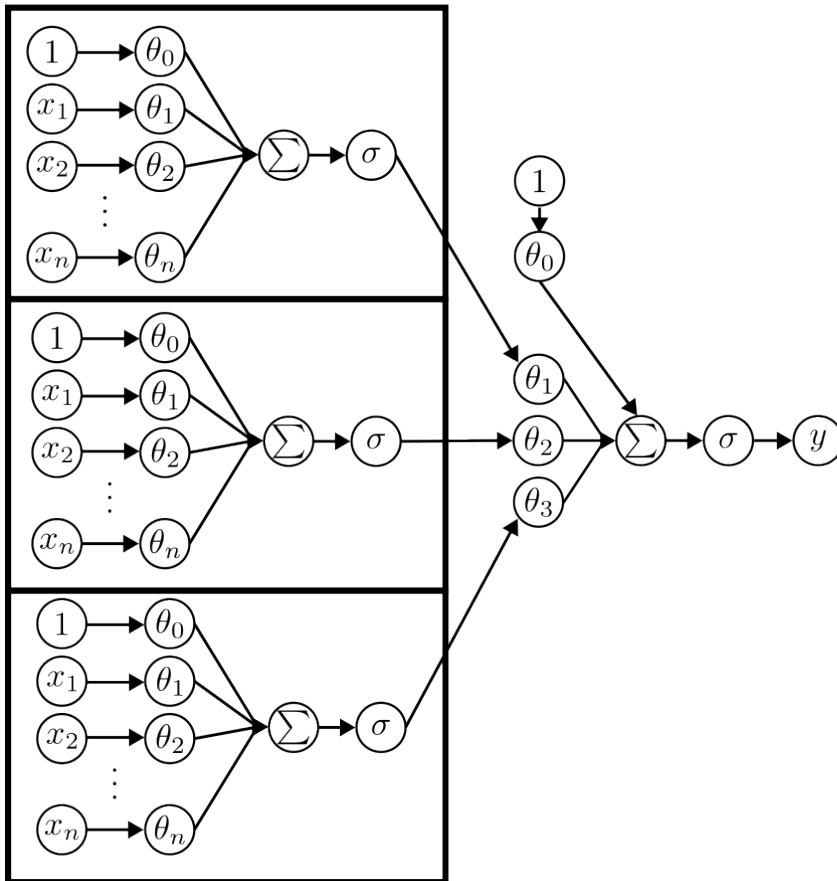
Adding neurons in **series** creates a **deep** neural network



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:

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

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

A single neuron:

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

Multiple neurons (wide):

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

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

A single neuron:

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

Multiple neurons (wide):

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

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:

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

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

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

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

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

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

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

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

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

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

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

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

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

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

A single neuron:

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

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

A single neuron:

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

Multiple neurons (deep):

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

A single neuron

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

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 $\boldsymbol{\theta}, \boldsymbol{\psi}, \boldsymbol{\rho}$

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

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

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

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

⋮

Written more plainly as

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

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

⋮

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

Written more plainly as

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

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

⋮

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

Implement XOR using a deep and wide neural network

Implement XOR using a deep and wide neural network

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

Implement XOR using a deep and wide neural network

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

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

Implement XOR using a deep and wide neural network

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

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

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?

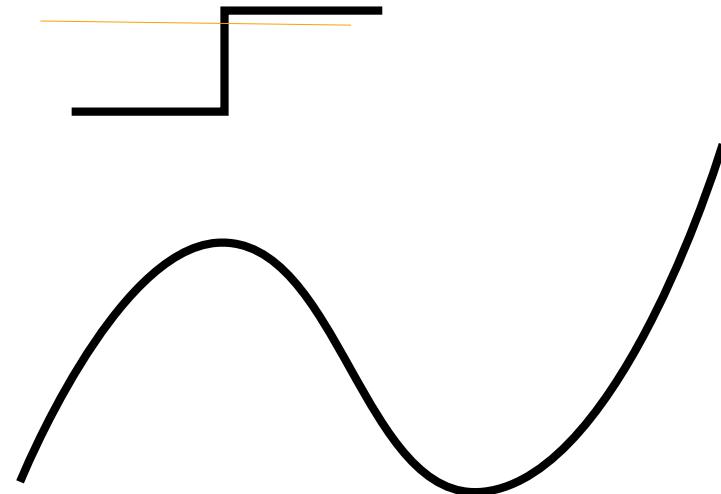
Consider a one-dimensional arbitrary function $g(x) = y$

We can approximate g using our neural network f

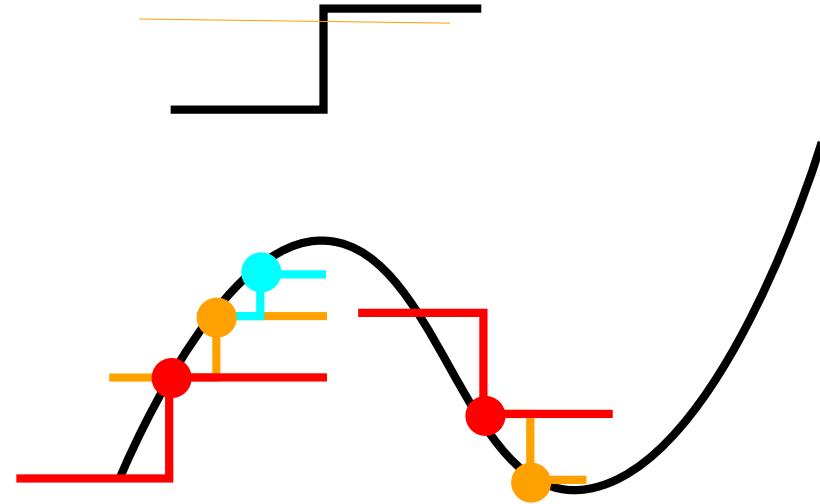
$$\begin{aligned} f(x_1, x_2, \theta) &= H(\theta_{3,0} \\ &\quad + \theta_{3,1} \cdot H(\theta_{1,0} + x_1 \theta_{1,1} + x_2 \theta_{1,2})) \\ &\quad + \theta_{3,2} \cdot H(\theta_{2,0} + x_1 \theta_{2,1} + x_2 \theta_{2,2})) \end{aligned}$$

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

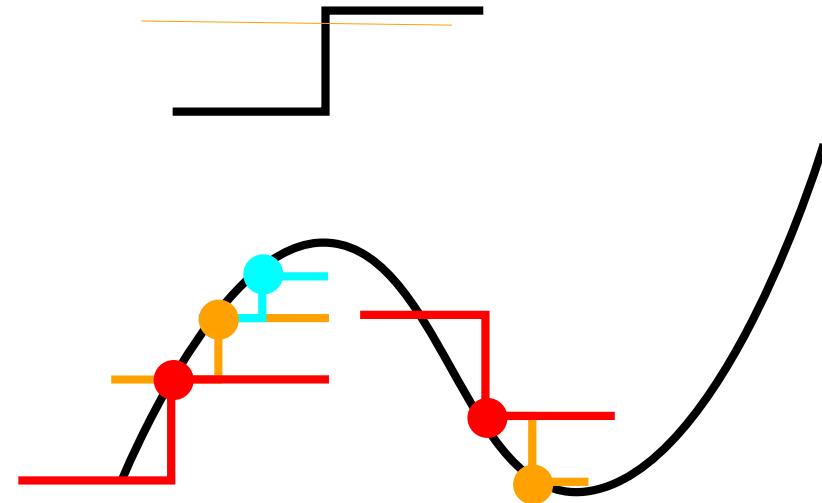
Proof Sketch: Approximate a function $g(x)$ using a linear combination of Heaviside functions



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



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

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 ε

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

It can approximate **any** continuous function to precision ε

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

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

It can approximate **any** continuous function to precision ε

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

As we increase the width and depth of the network, ε shrinks

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

It can approximate **any** continuous function to precision ε

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

As we increase the width and depth of the network, ε shrinks

$$g\left(\begin{array}{c} \text{Dog's face} \end{array}\right) = \text{Dog}$$

$$g\left(\begin{array}{c} \text{Muffin} \end{array}\right) = \text{Muffin}$$

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

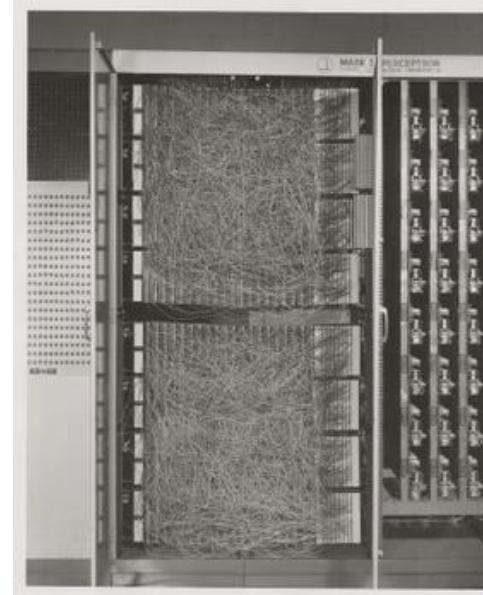
All the models we examine in this course will use this neural network structure internally

- Transformers
- Graph neural networks

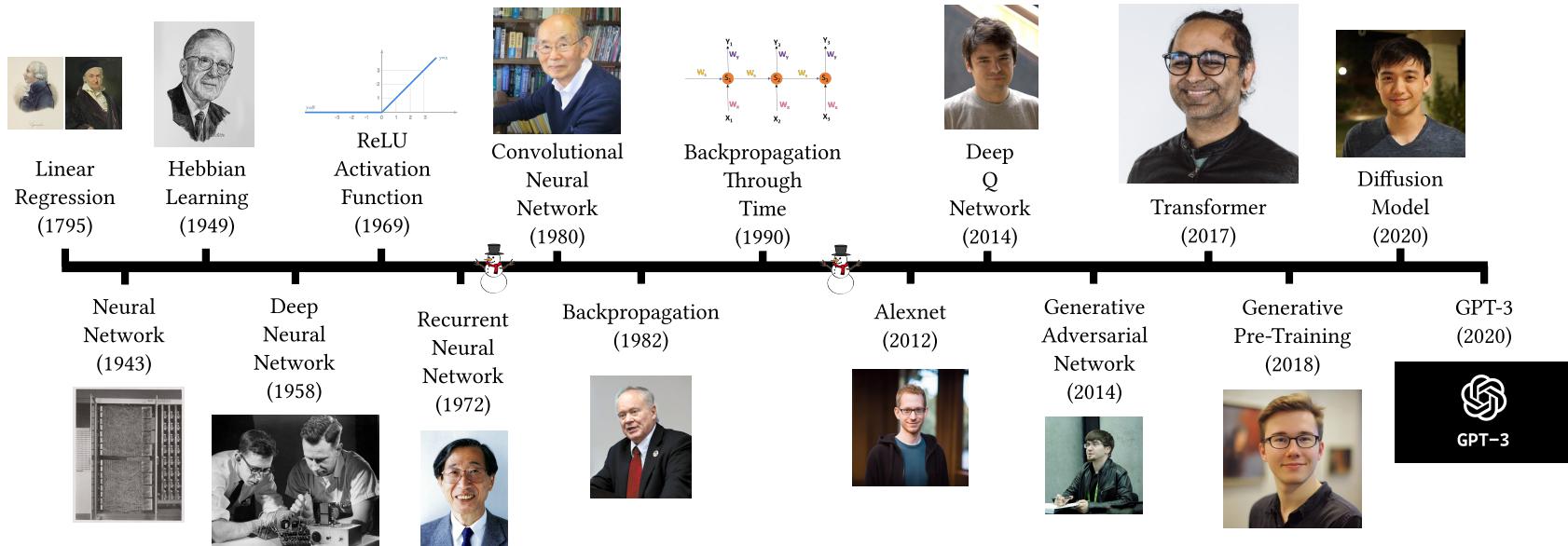
Relax

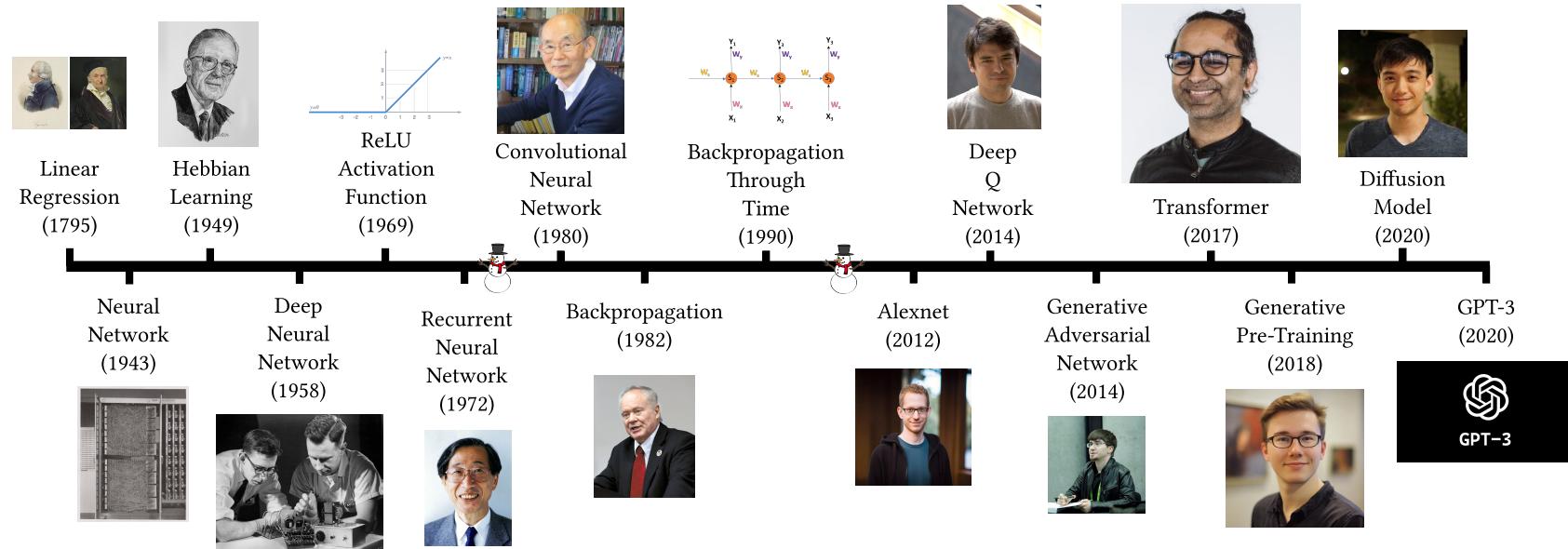
We call this form of a neural network a **feedforward network** or **perceptron** (invented in 1943)

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



20×20 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 θ

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

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

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

Here are the equations for one neural network layer

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

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

Let us do this in colab! https://colab.research.google.com/drive/1bLtf3QY-yROIif_EoQSU1WS7svd0q8j7?usp=sharing