

Introduction to Machine Learning: Neural Networks

Stefania Sarno

January 6, 2021

Outline

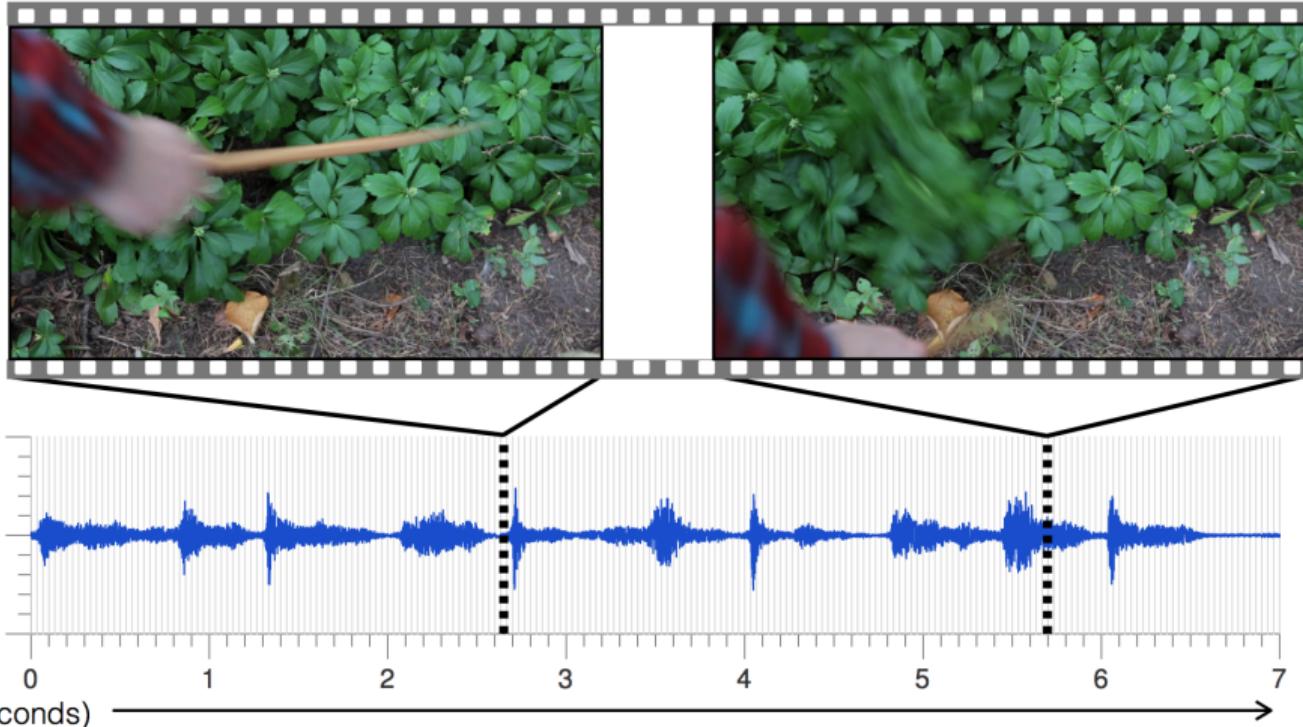
1. Recent successes of neural networks and deep learning
2. Historical development of neural networks
3. Perceptron (learning rule and limitations)
4. Multi-layer Perceptrons and the Universal Approximation Theorem
5. Backpropagation Algorithm

Automatic Translation

- Breakthrough of Google Translate (2016)
- But neural network can also automatically translate from images



Sound Generation from Silent Videos



<https://news.mit.edu/2016/artificial-intelligence-produces-realistic-sounds-0613>

Image segmentation

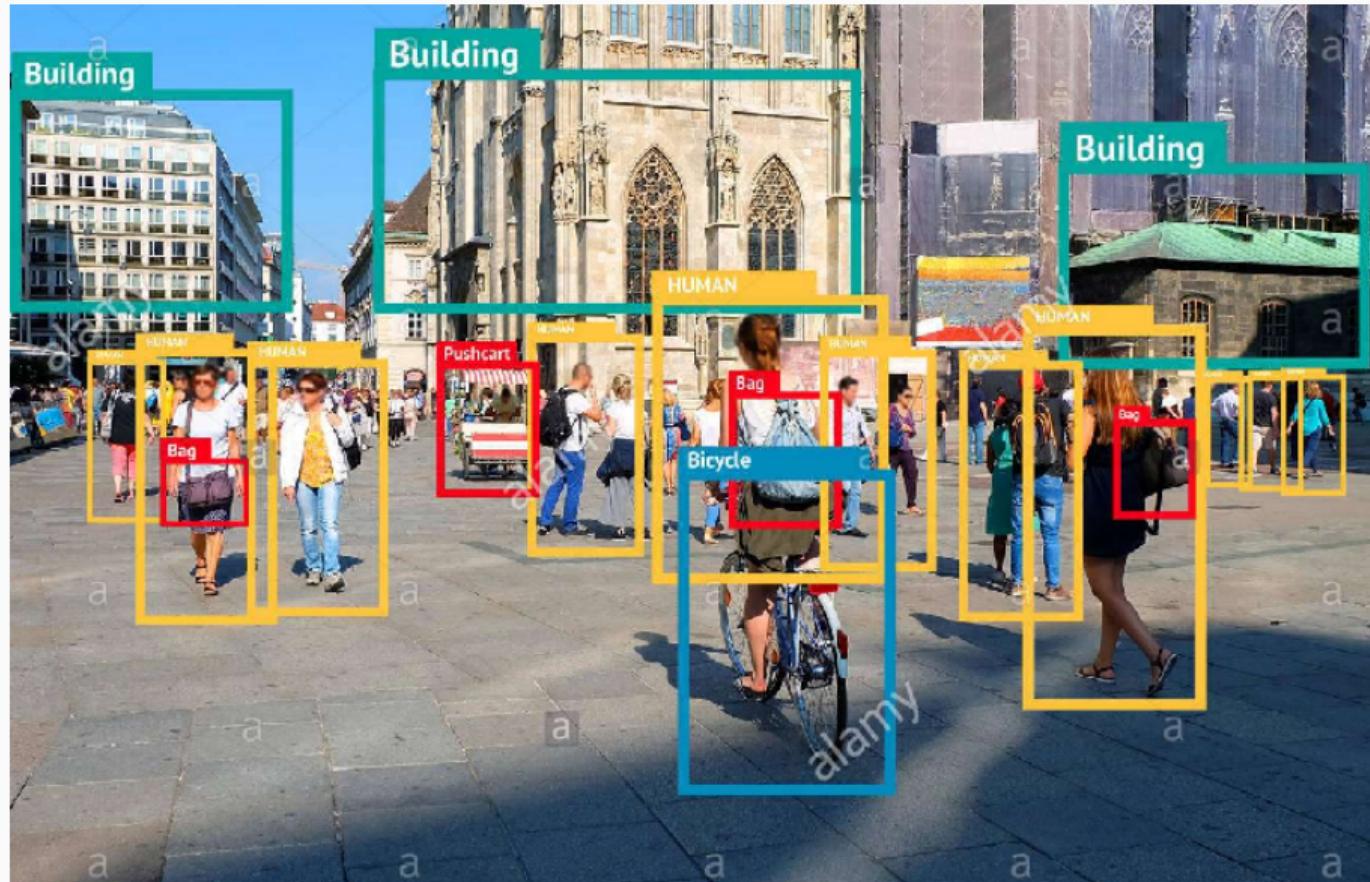
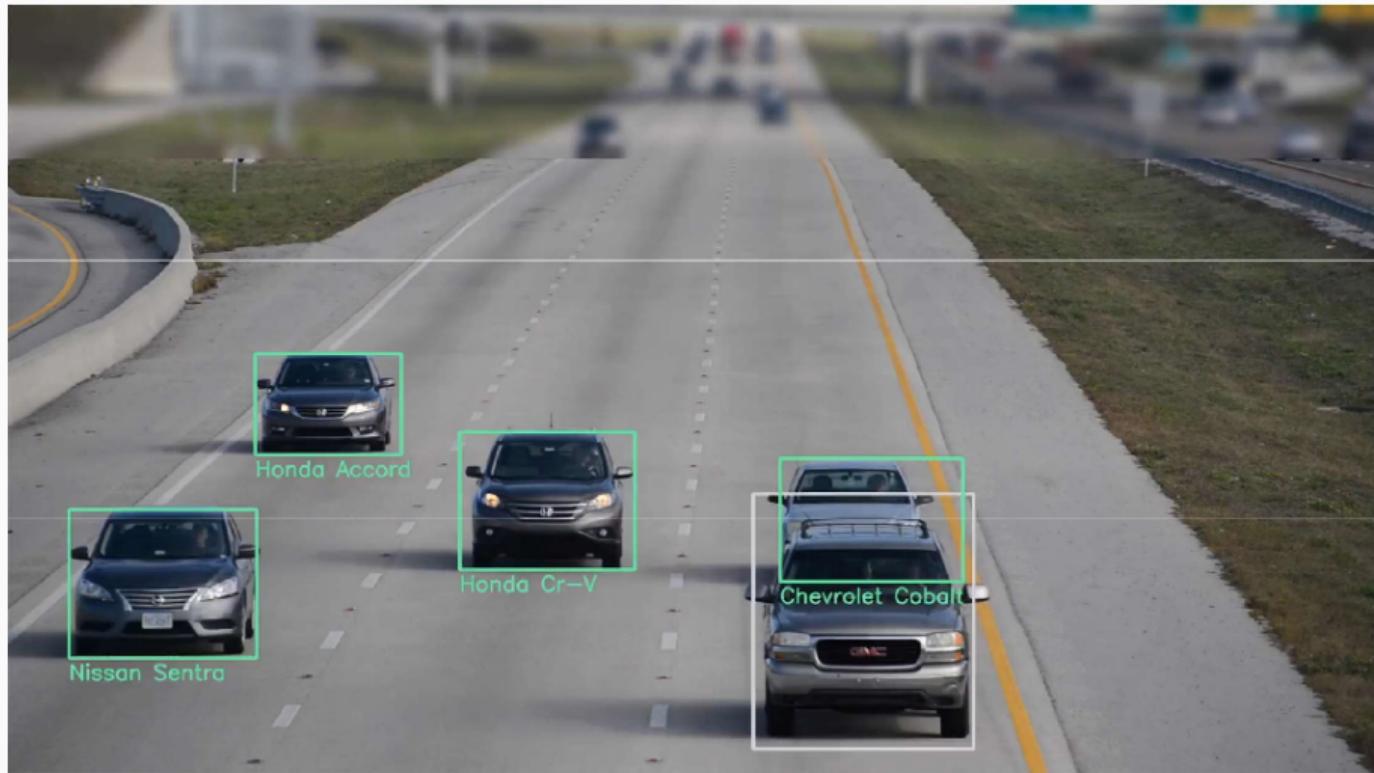


Image Recognition

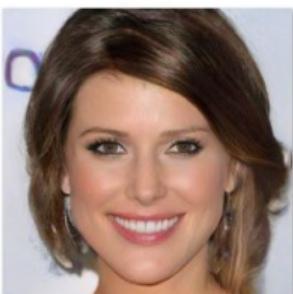
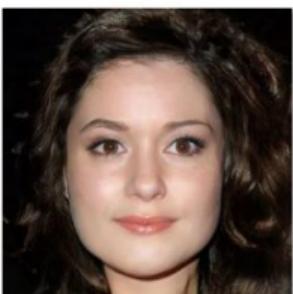
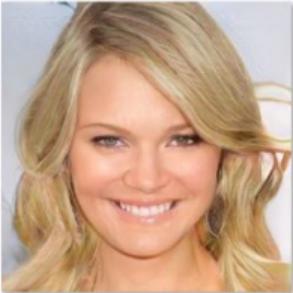


<https://www.sighthound.com/technology/>

9

5

Fake Face Generation



<https://medium.com/datadriveninvestor/artificial-intelligence-gans-can-create-fake-celebrity-faces-44fe80d419f7>

Image Captions Generation



"man in black shirt is playing guitar."



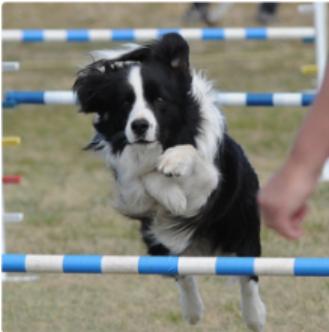
"construction worker in orange safety vest is working on road."



"two young girls are playing with lego toy."



"girl in pink dress is jumping in air."

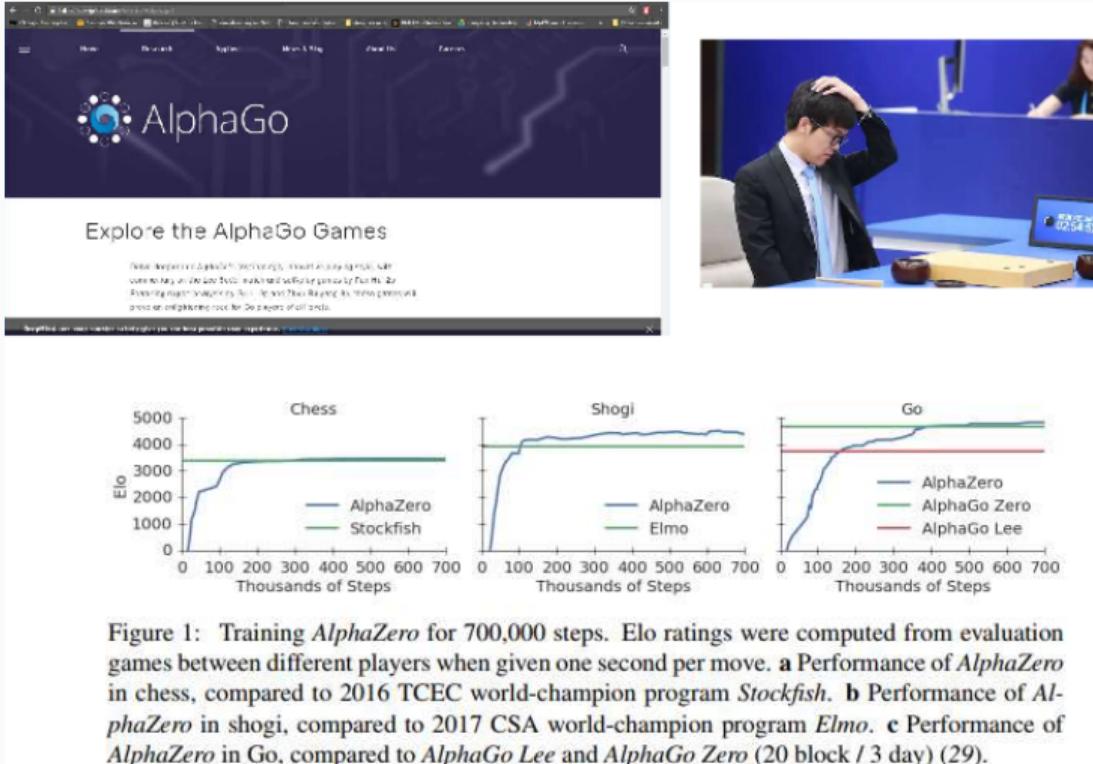


"black and white dog jumps over bar."



"young girl in pink shirt is swinging on swing."

Alpha Go

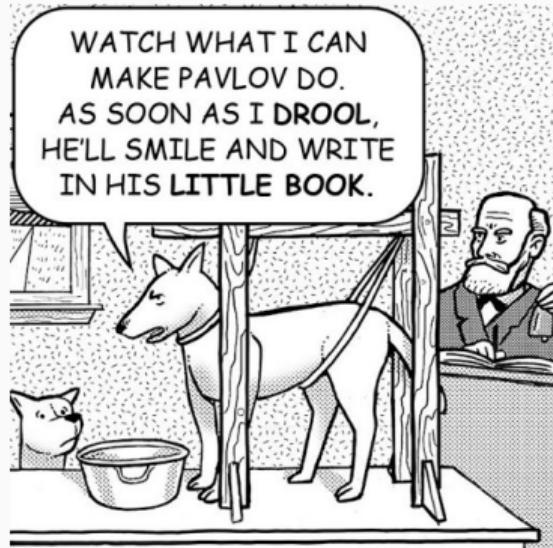


What are Neural Networks?

Everything began with the study of the brain



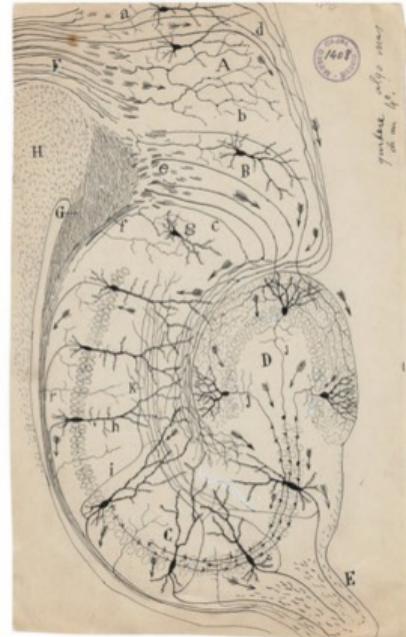
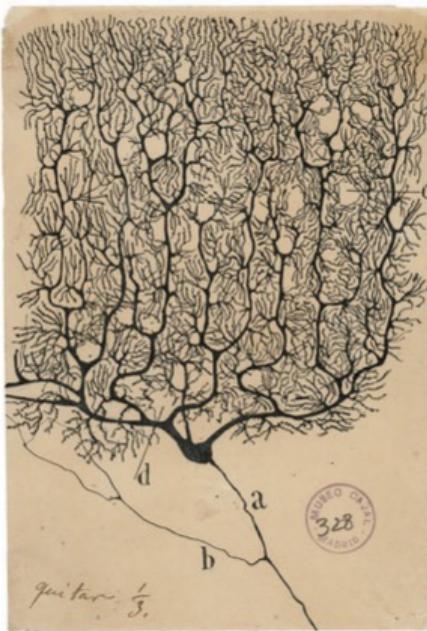
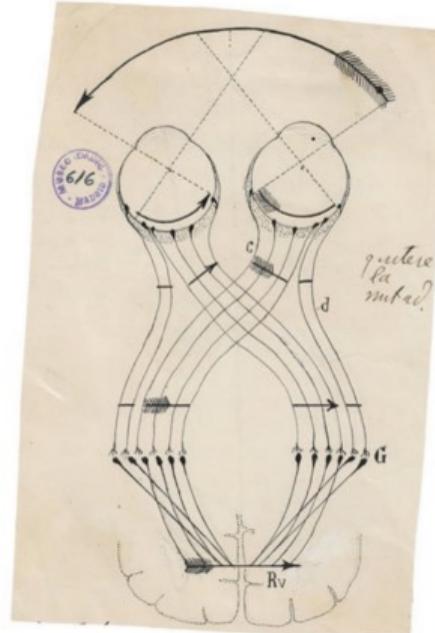
Earliest Model of Cognition: Associationism



- The ability to associate ideas is the only mental process
- We learn by cause and effect, contiguity, resemblance
- Behaviour is learned from repeated associations of actions with feedback

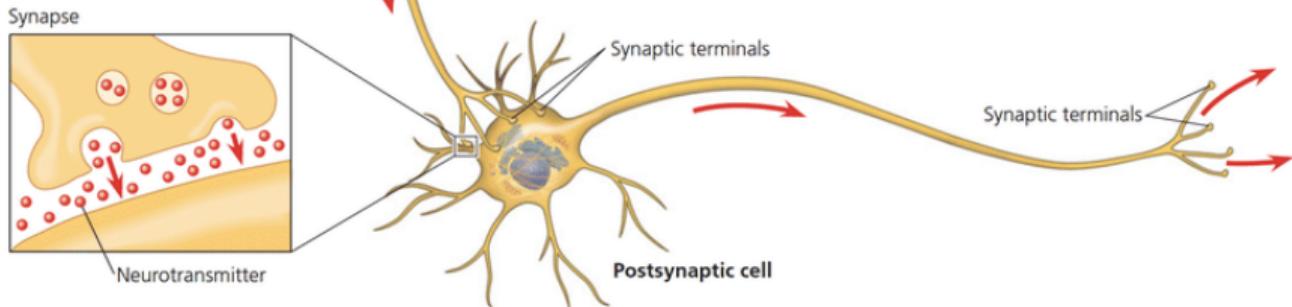
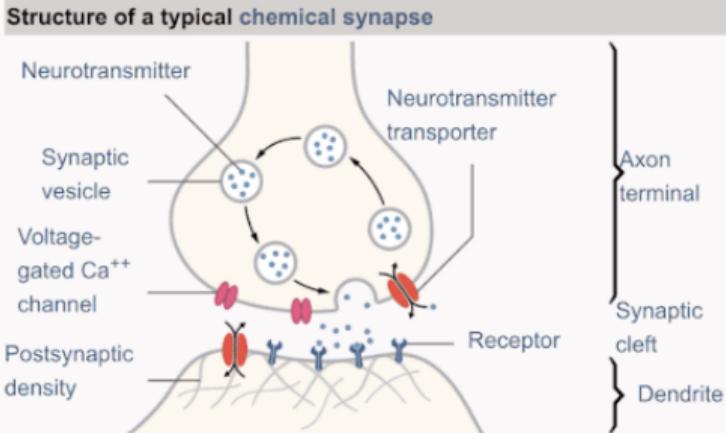
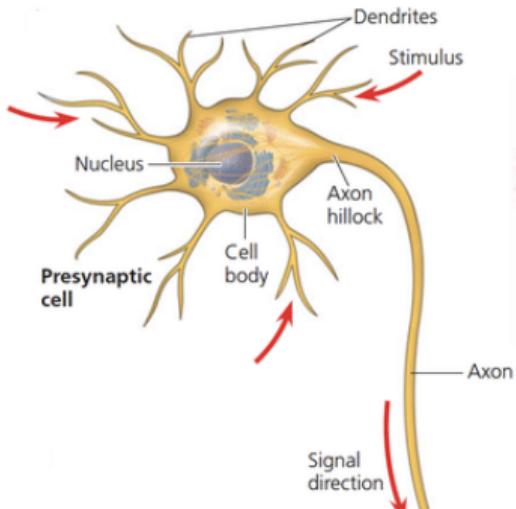
But where are associations stored?

The Discovery of the Neurons (late 1800)

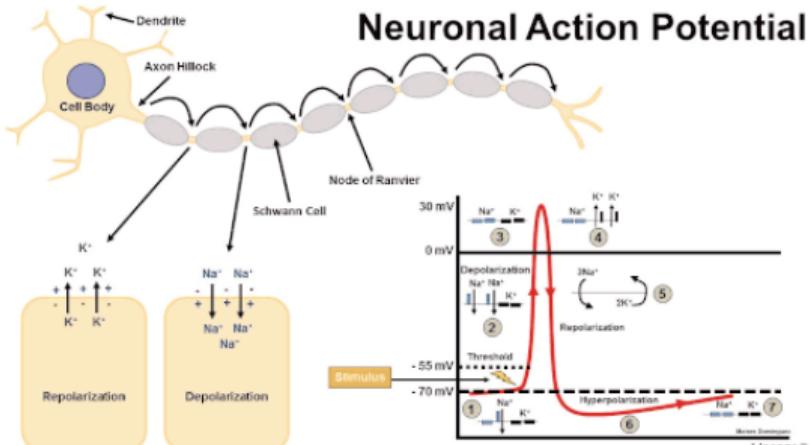
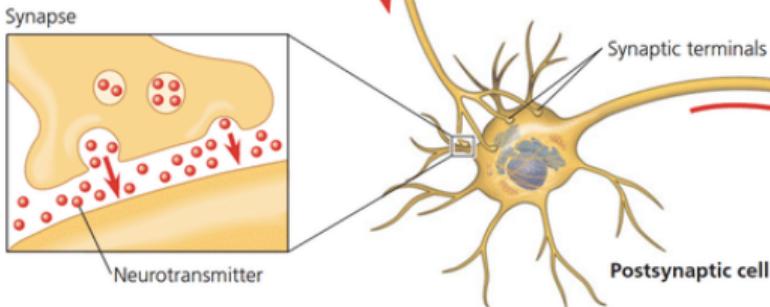
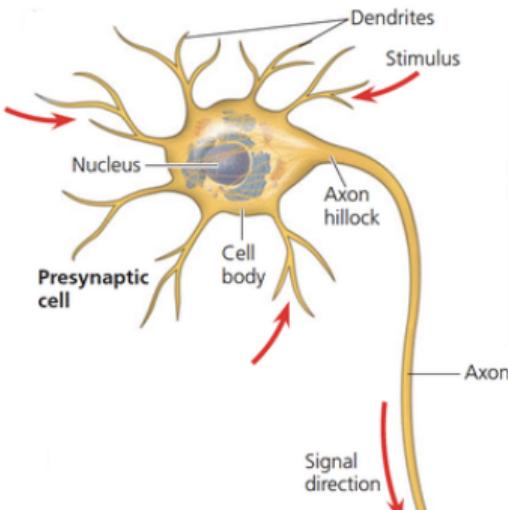


Illustrations by Santiago Ramón y Cajal, the Spanish neuroscientist, from the book 'The Beautiful Brain.' From left: A diagram suggesting how the eyes might transmit a unified picture of the world to the brain; a purkinje neuron from the human cerebellum; and a diagram showing the flow of information through the hippocampus in the brain.

How Real Neurons Work

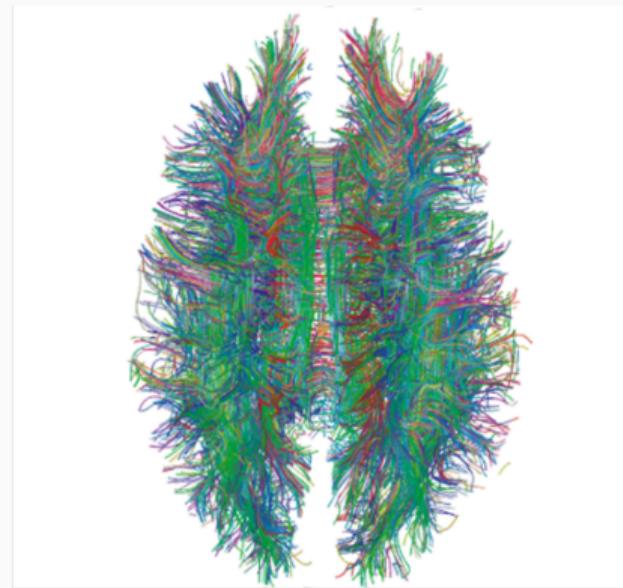


The Action Potential



The Brain as a Connectionist Machine

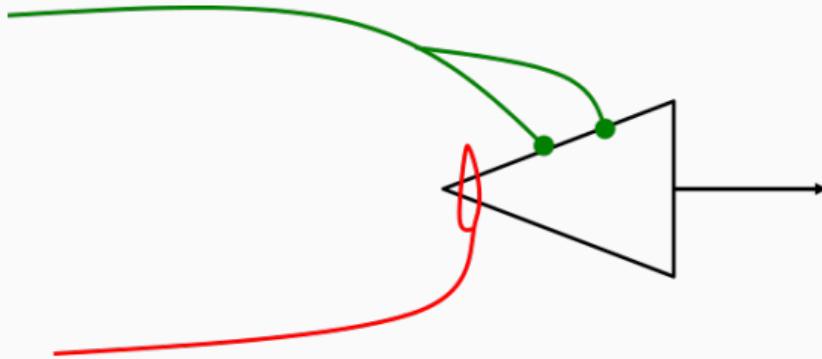
- There are ~100 billion neurons
- Each neuron has an average of 7000 synaptic connections
- 10^{15} synapses
- 1.3 – 1.4 kg



Gigandet X, Hagmann P, Kurant M, Cammoun L, Meuli R, et al. *PLoS ONE* (2008)

Connectionism: Everything is in the connections

McCulloch and Pitts model (1943)



"*A Logical Calculus of the Ideas Immanent in Nervous Activity*", Bulletin of Mathematical Biophysics

- **Excitatory synapse:** Transmit weighted input to the neuron
- **Inhibitory synapse:** If receiving signal, it prevents the neuron from firing

Hebbian Learning

Donald Hebb: *Organization of Behaviour*, 1949:

Let us assume that the persistence or repetition of a reverberatory activity (or "trace") tends to induce lasting cellular changes that add to its stability. ... When an axon of cell A is near enough to excite a cell B and repeatedly or persistently takes part in firing it, some growth process or metabolic change takes place in one or both cells such that A's efficiency, as one of the cells firing B, is increased.

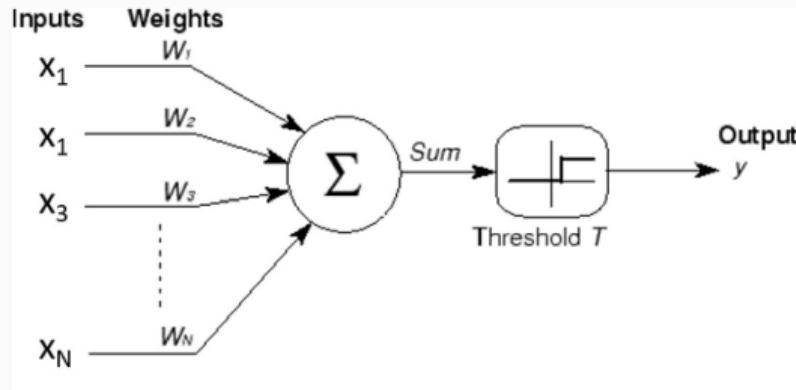
Neurons that fire together wire together

If neuron x repeatedly triggers neuron y their connection strength will change as:

$$w_{xy} = w_{xy} + \alpha xy$$

The Perceptron: simplified model and learning algorithm

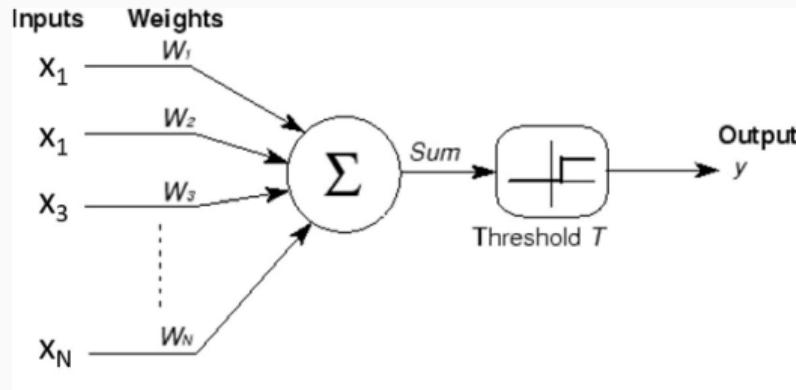
Frank Rosenblatt, 1958 (Boolean tasks and proved convergence)



$$y = \begin{cases} 1 & \text{if } \sum_i w_i x_i - T \geq 0 \\ 0 & \text{else} \end{cases}$$

The Perceptron: simplified model and learning algorithm

Frank Rosenblatt, 1958 (Boolean tasks and proved convergence)

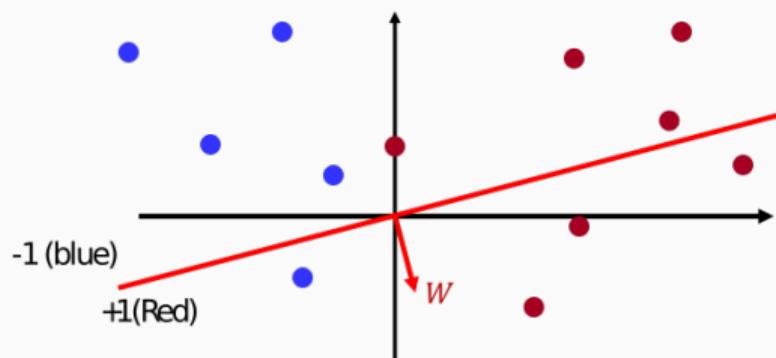


$$y = \begin{cases} 1 & \text{if } \sum_i w_i x_i - T \geq 0 \\ 0 & \text{else} \end{cases}$$

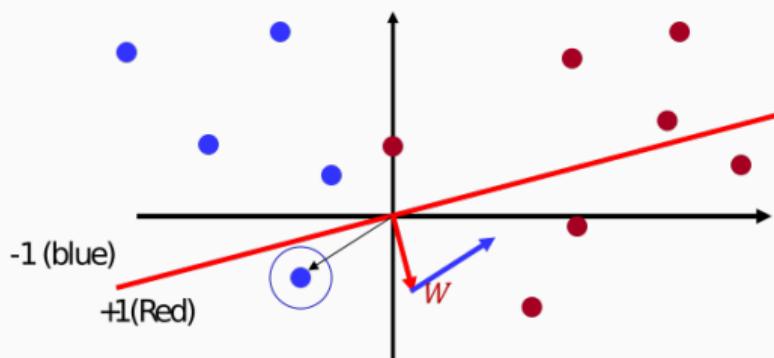
- Update the weights whenever the perceptron output is wrong as:

$$\mathbf{w} = \mathbf{w} + \alpha [d(\mathbf{x}) - y(\mathbf{x})] \mathbf{x}$$

The Perceptron: the Learning Algorithm

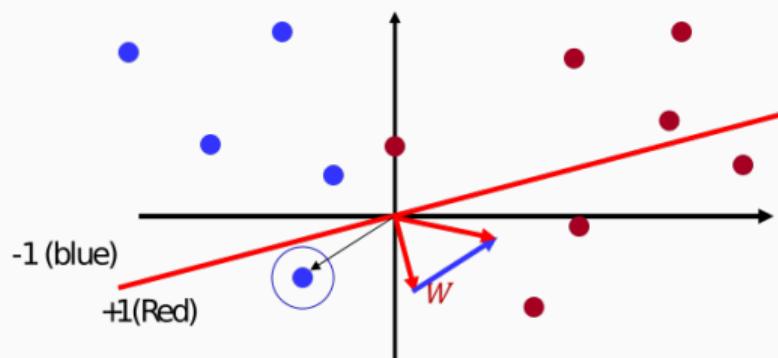


The Perceptron: the Learning Algorithm

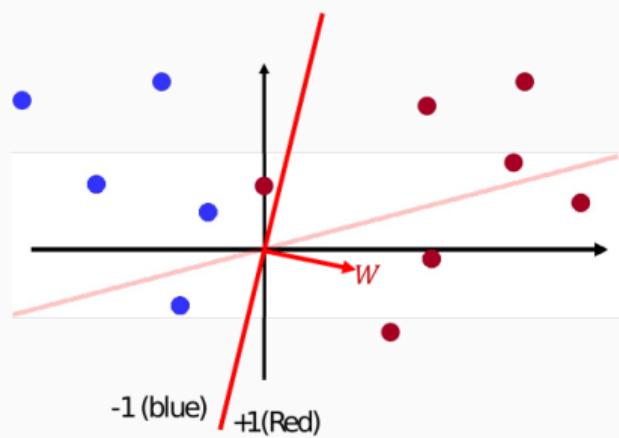


Misclassified negative instance, subtract it from W

The Perceptron: the Learning Algorithm

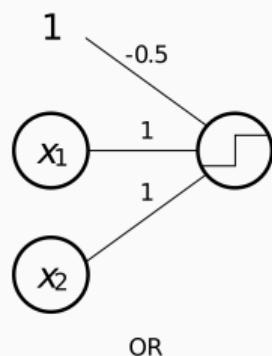


The Perceptron: the Learning Algorithm

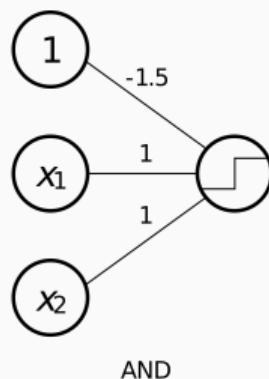


The Boolean Gates and the XOR Problem

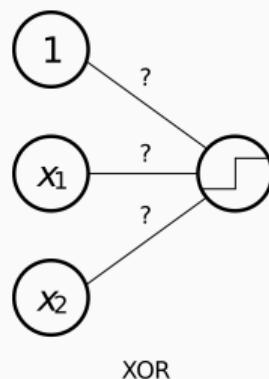
Minsky and Papert, 1968



OR



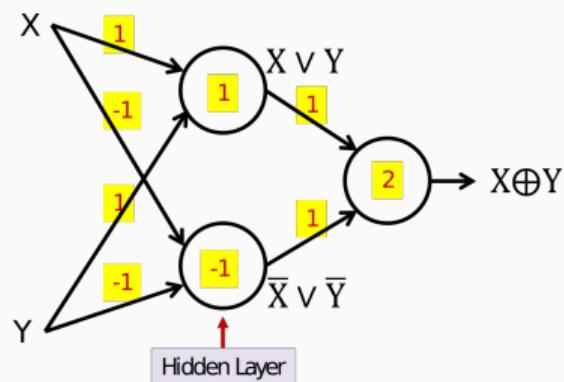
AND



XOR

Solving the XOR Problem: the Multilayer Perceptron (MLP)

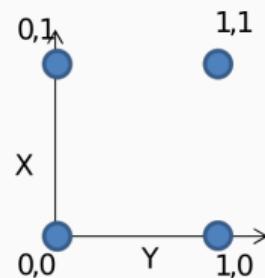
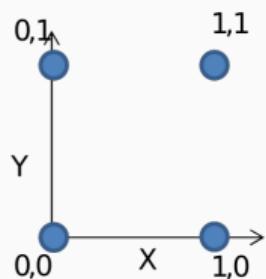
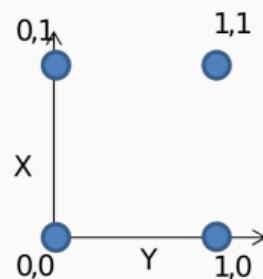
Minsky and Papert, 1968



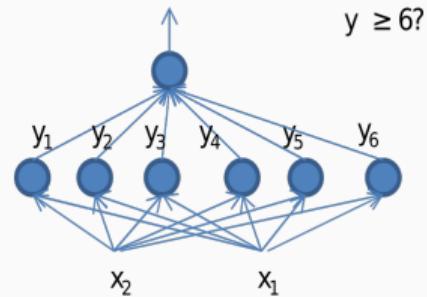
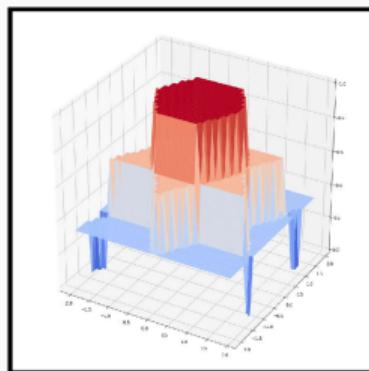
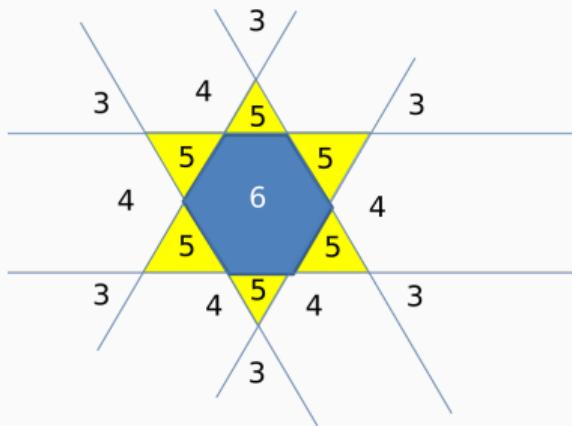
The XOR Problem

Minsky and Papert, 1968

Why do we have this?

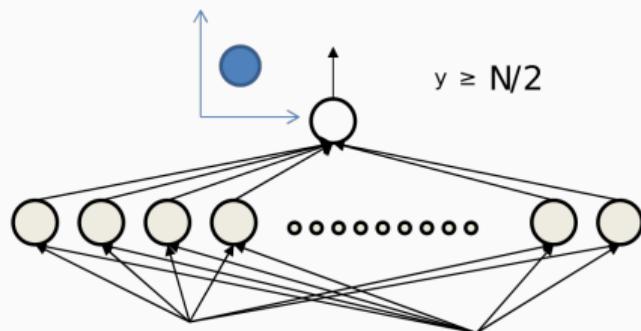
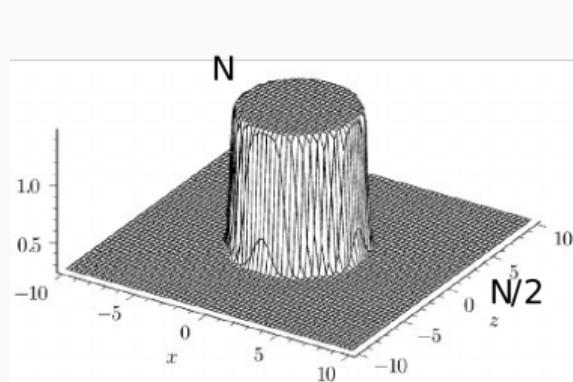


MLP as Function Approximators



MLP as Function Approximators

For a very large number N of hidden units

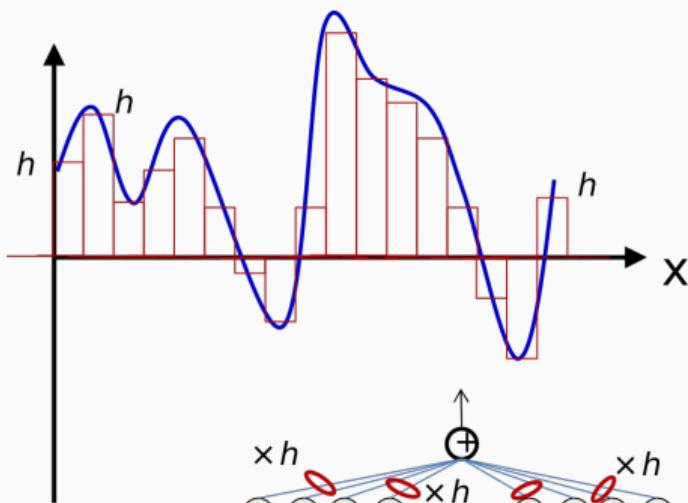
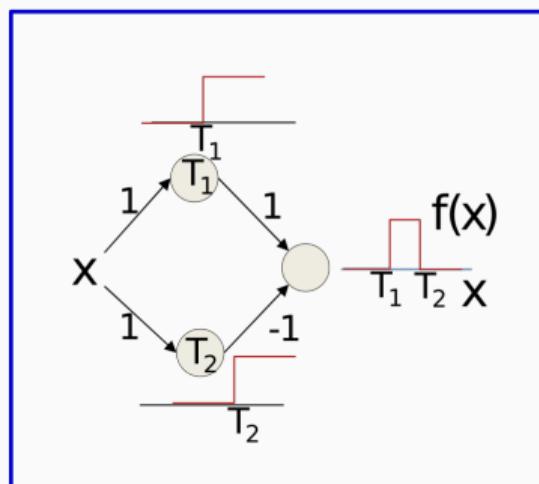


MLP as Function Approximators

Once I have a circle I can do any weird boundary

MLP as Function Approximators

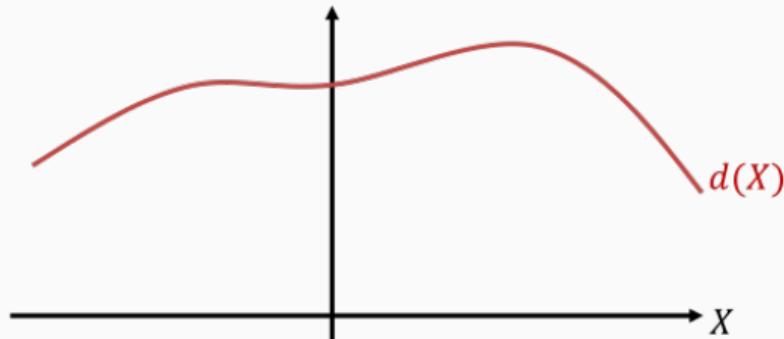
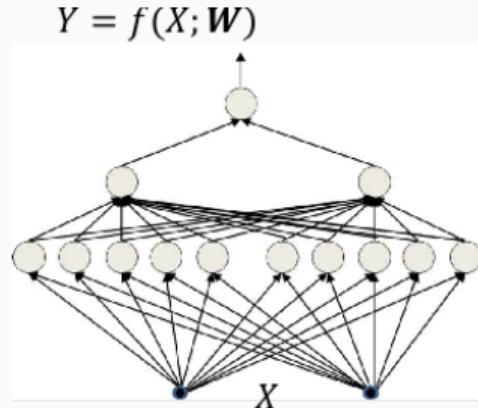
We can also approximate real function with arbitrary precision



Universal Approximation Theorem

- A feed-forward network with a single hidden layer containing a finite number of neurons can approximate **continuous** functions on compact subsets of \mathbb{R}^n , under mild assumptions on the activation function.
- Proved by George Cybenko in 1989 for logistic activation function

Neural Networks as Universal Approximators

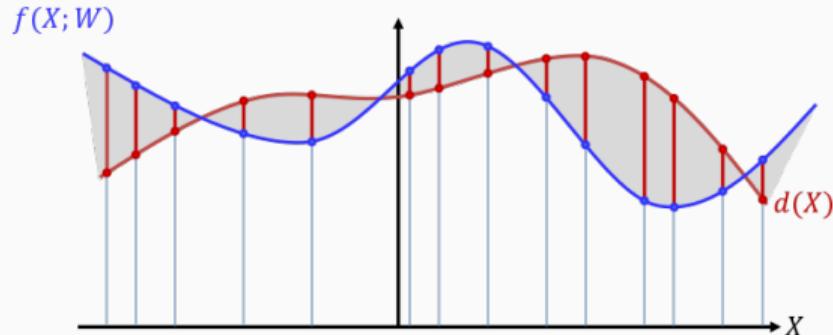


- Neural networks can approximate any function
- We want to minimize the "distance" between the function produced by our network and the function we want to mimic

$$\text{distance}[f(X, \mathbf{W}), d(X)]$$

- We have to determine the weights and bias to get a good approximation

The Sampling Idea



Problem: we do not know $d(X)$ at all points X

Sampling Solution:

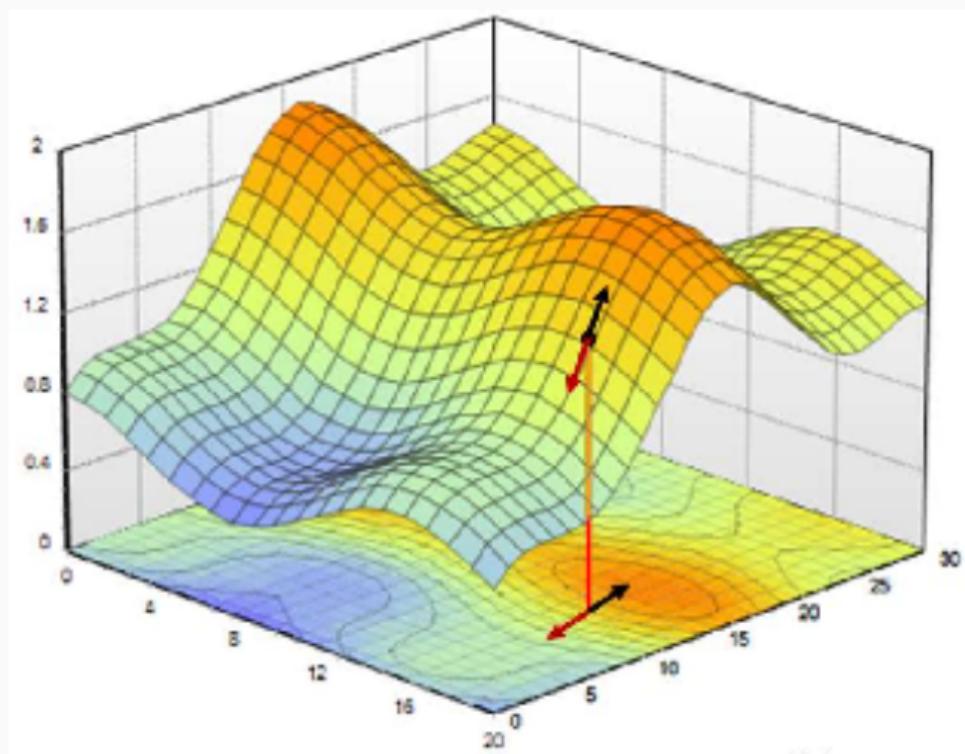
- Provided enough "training" samples we the distance can be approximated as:

$$Cost(\mathbf{W}) = \frac{1}{N} \sum_i \text{distance}[f(X_i, \mathbf{W}), d(X_i)]$$

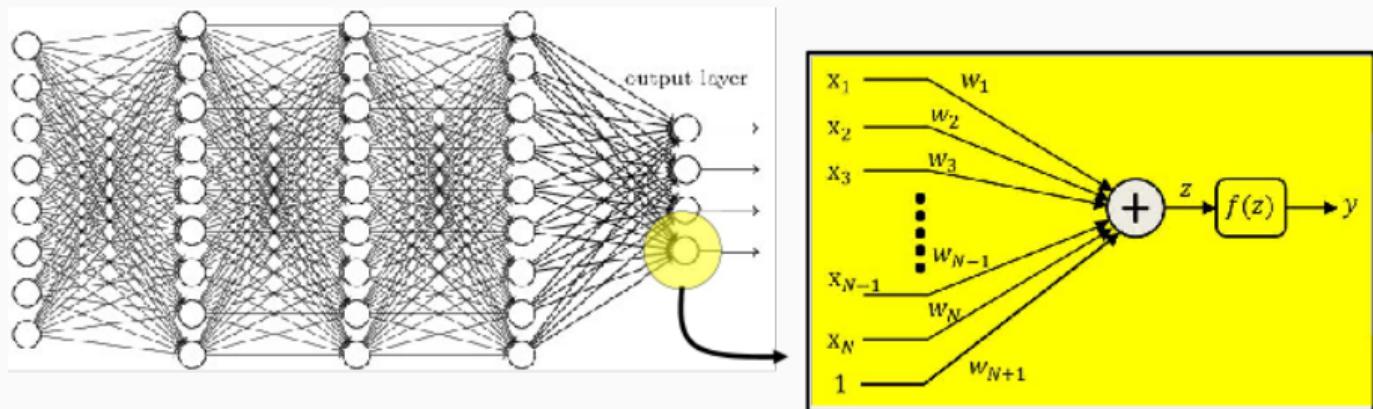
- The weights can be estimated by minimizing the cost:

$$\tilde{\mathbf{W}} = \operatorname{argmin}_{\mathbf{W}} Cost(\mathbf{W})$$

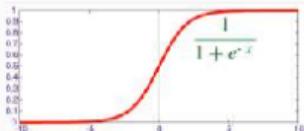
Finding the Minimum by Gradient Descend



Activation Functions for Multilayer Perceptrons



Activation Functions for Multilayer Perceptrons



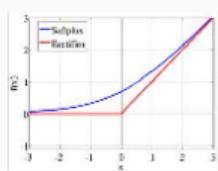
$$f(z) = \frac{1}{1 + \exp(-z)}$$

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



$$f(z) = \tanh(z)$$

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



$$f(z) = \begin{cases} z, & z \geq 0 \\ 0, & z < 0 \end{cases}$$

$$f'(z) = \begin{cases} 1, & z \geq 0 \\ 0, & z < 0 \end{cases}$$

$$f(z) = \log(1 + \exp(z))$$

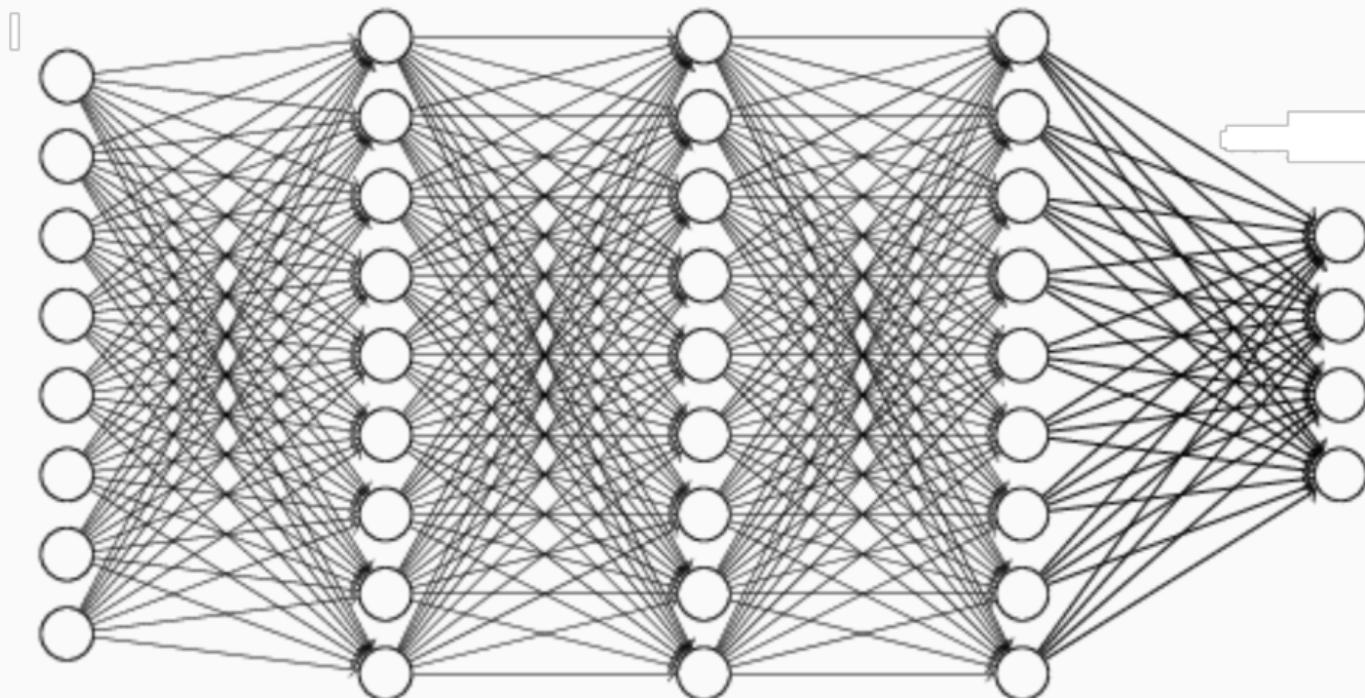
$$f'(z) = \frac{1}{1 + \exp(-z)}$$

The Cost Function of our Minimization Problem

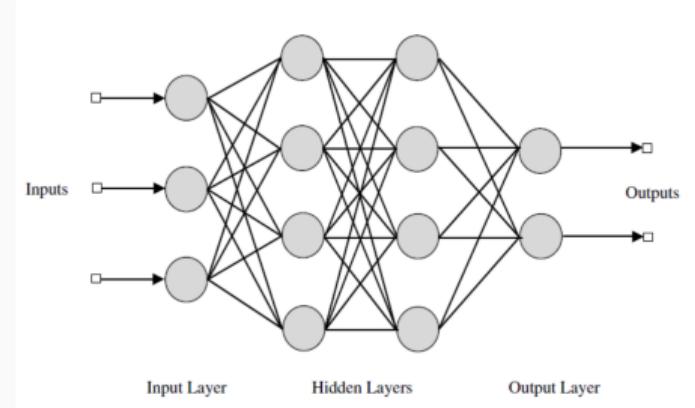
We already did this for:

1. Linear Regression
2. Logistic Regression

Fully-connected Feedforward Networks



The Backpropagation Algorithm 1



Summary of the forward pass

$$o^{(1)} = x$$

$$p^{(2)} = \Theta^{(1)} o^{(1)}$$

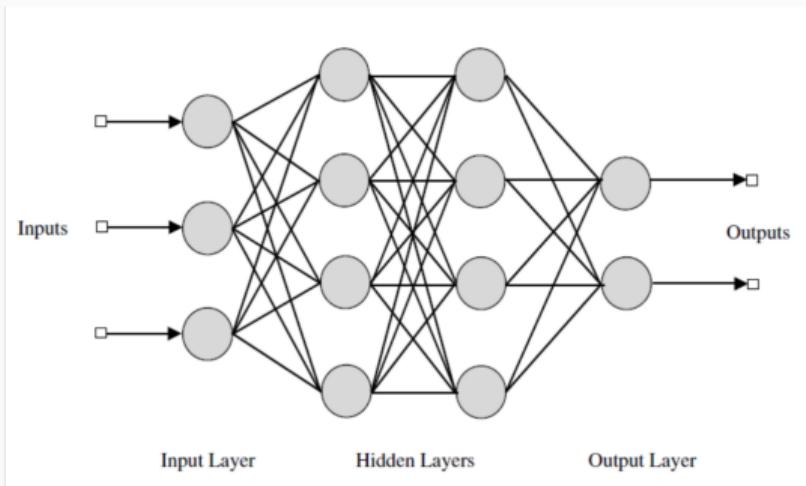
$$o^{(2)} = g(p^{(2)})$$

$$p^{(3)} = \Theta^{(2)} o^{(2)}$$

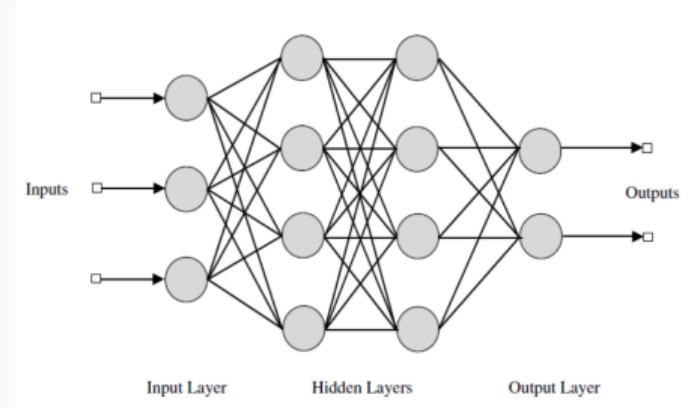
$$o^{(3)} = g(p^{(3)})$$

$$p^{(4)} = \Theta^{(3)} o^{(3)}$$

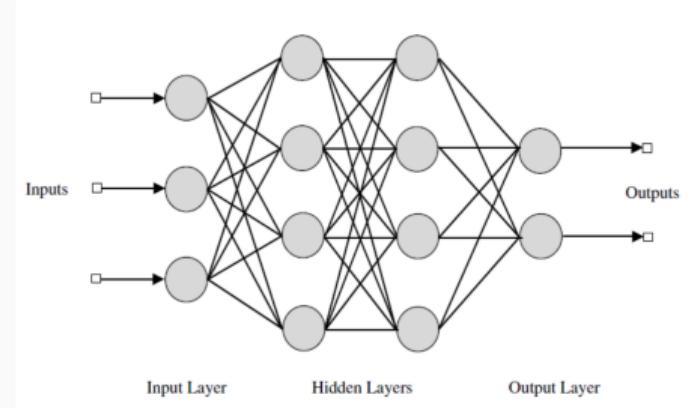
$$o^{(4)} = g(p^{(4)}) = h_{\theta}(x)$$



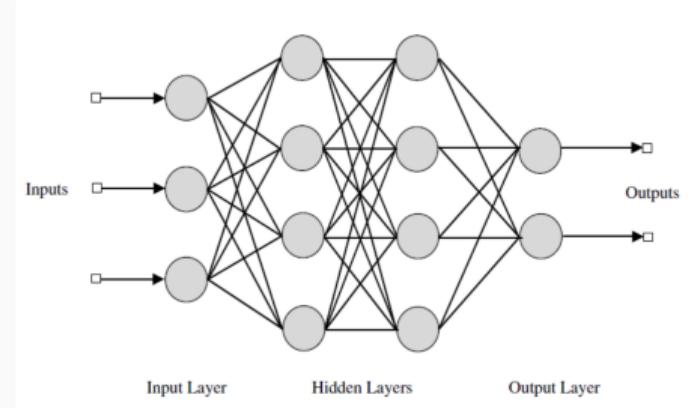
The Backpropagation Algorithm 2



The Backpropagation Algorithm 3



The Backpropagation Algorithm 4



The Backpropagation Algorithm 5

Summary of the backward pass

- Error: $E = \frac{1}{2} (h_\theta(x) - y)^2$

- Vectors δ 's

$$\delta^{(L)} = (h_\theta - y) g(p^{(L)})$$

$$\delta^{(l)} = (\Theta^{(l)})^T \delta^{(l+1)} \cdot g(p^{(l)})$$

- Gradient of error: $\frac{\partial E}{\partial \Theta_{ij}^{(l)}} = \delta_i^{(l+1)} o_j^{(l)}$

About the Gradient Descent

- Gradient Descent
- Stochastic Gradient Descent
- Mini-batch gradient descent