

Introduction to Machine Learning

SCP8084699 - LT Informatica



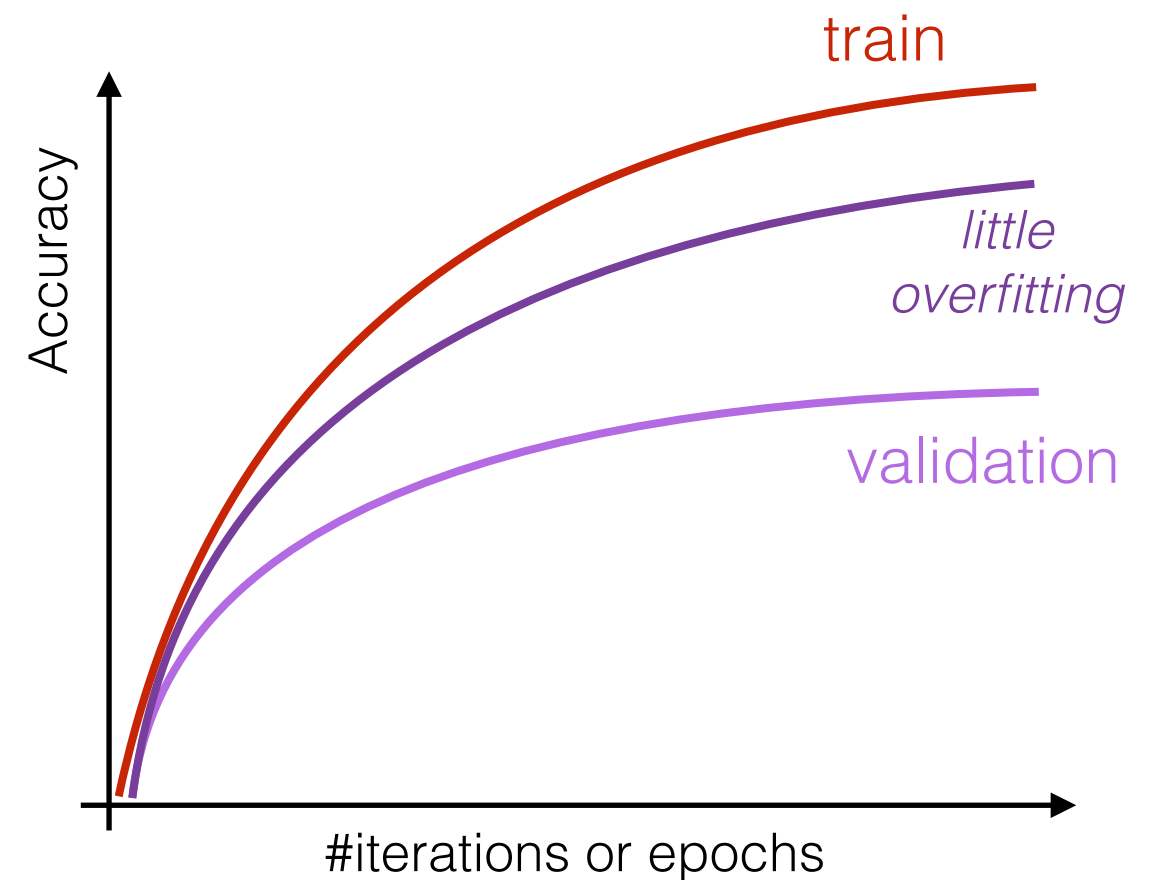
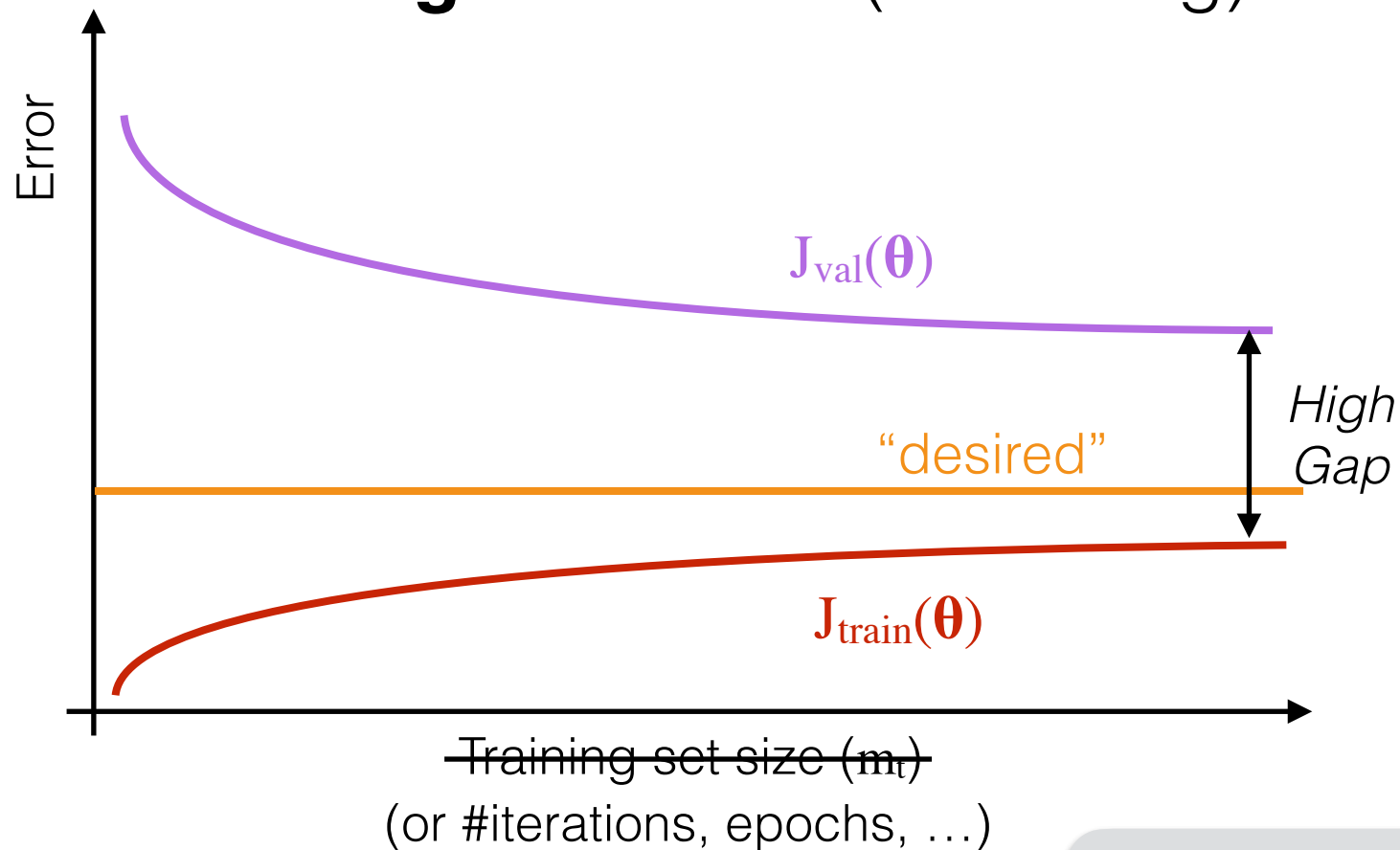
Evaluation, Learning Curves & Babysitting

Prof. Lamberto Ballan

Recap: Learning Curves

- You can compute learning curves w.r.t. different “dimensions” (e.g. evaluation measures, no. samples)

High variance (overfitting)

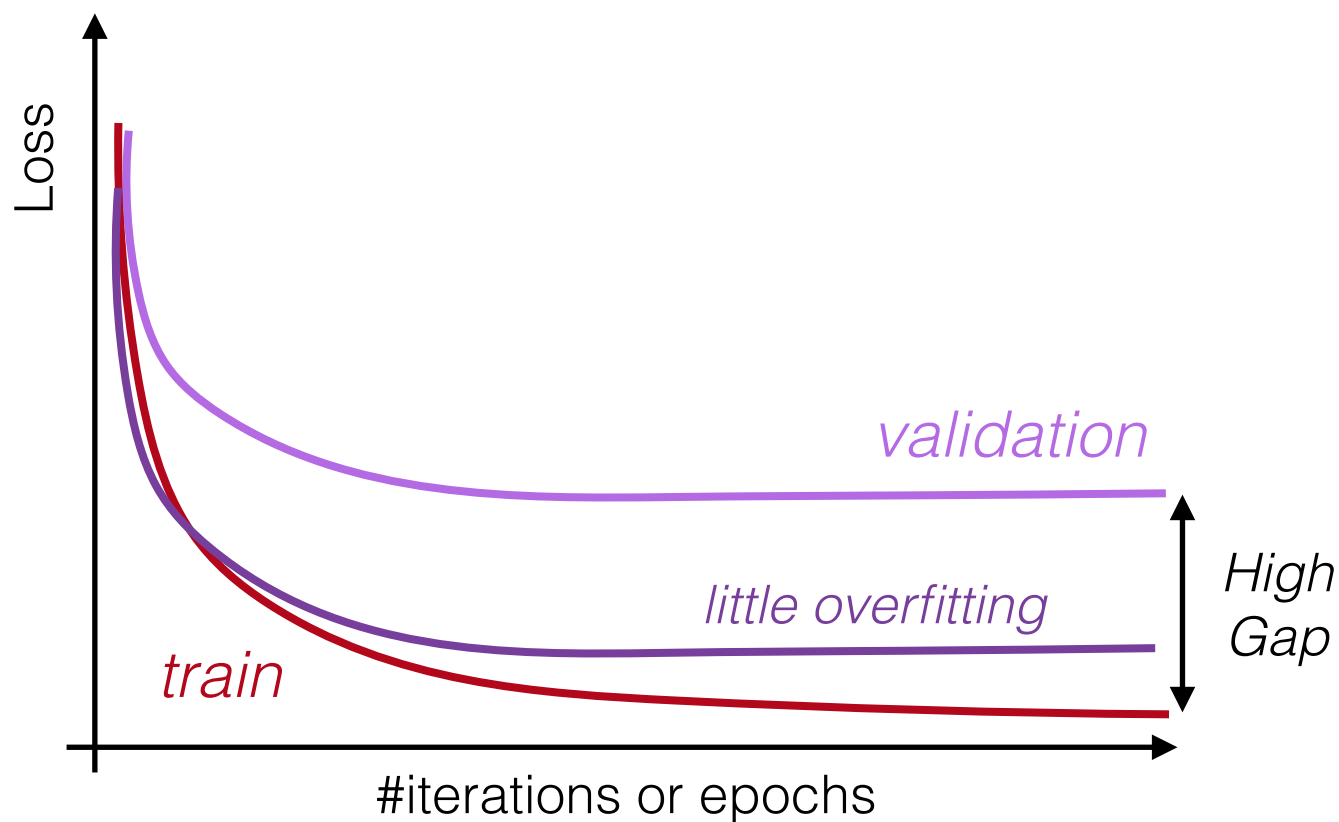


epoch = a pass of the entire set of examples

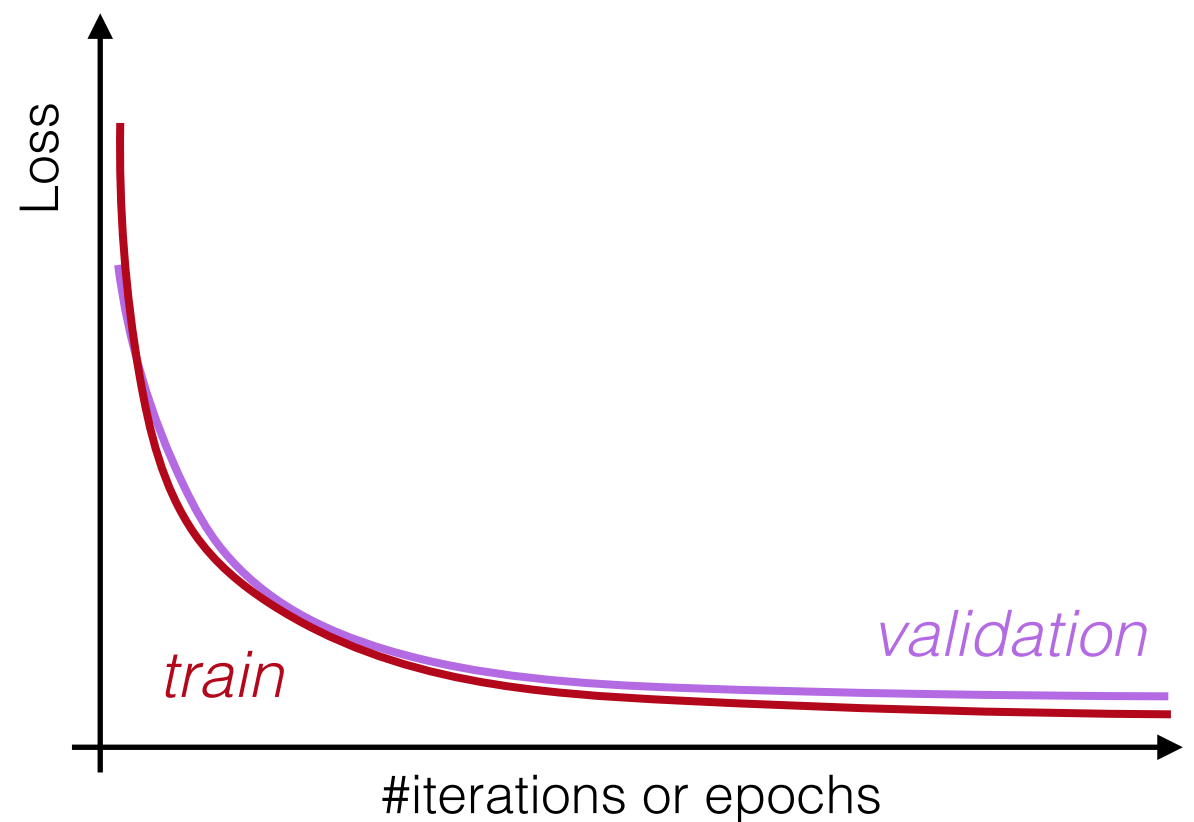
Learning Curves

- Often learning curves are plotted w.r.t. the loss

(overfitting)



(great fit)



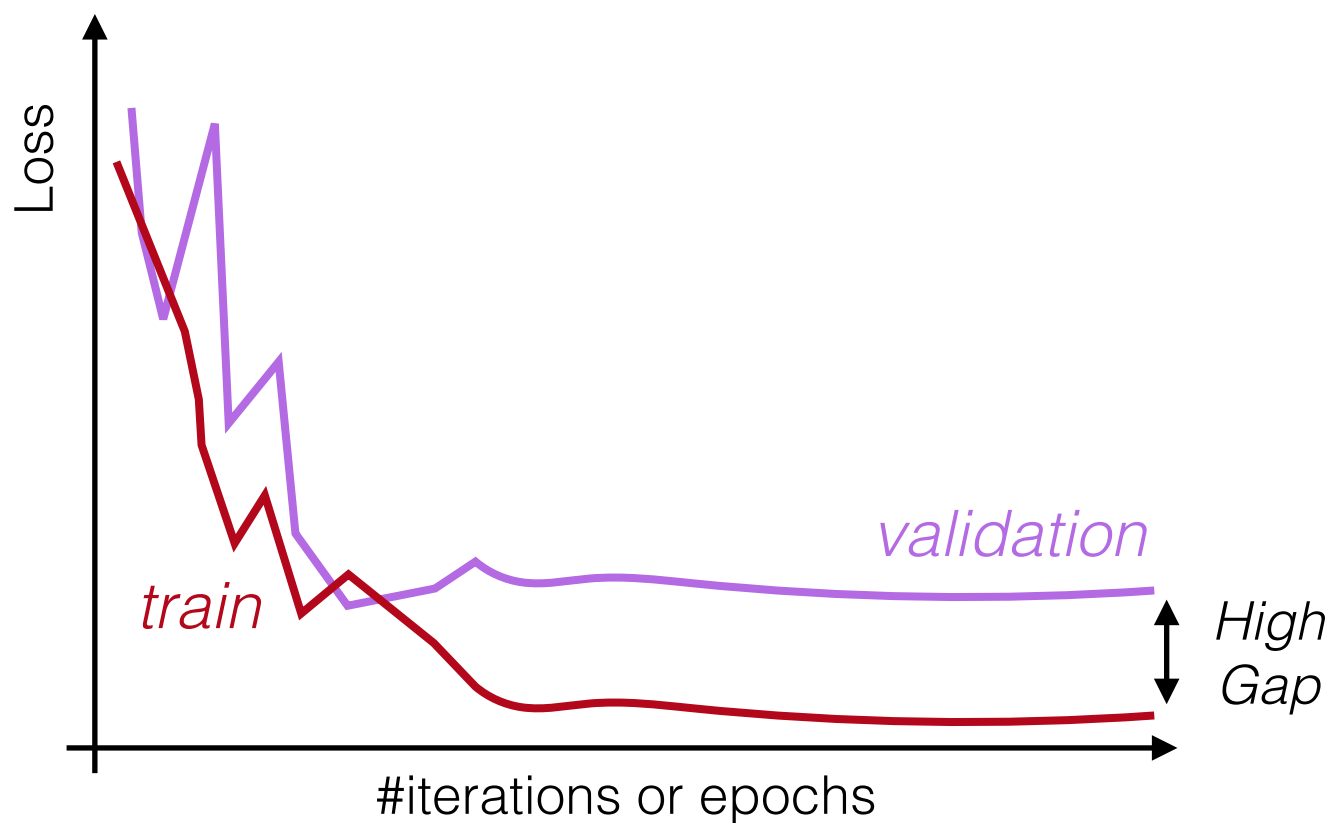
What to do next

- Debugging (and babysitting) a learning algorithm:
 - Suppose you have implemented a regularized linear regression model for predicting housing prices
 - It doesn't work on new data; what should you do next?
 - You can get more training data → *Fixes high variance*
 - Try smaller set of features → *Fixes high variance*
 - Try getting more features → *Fixes high bias*
 - Try adding complexity to the model (e.g. polynomial features)
 - Try decreasing λ → *Fixes high bias*
 - Try increasing λ → *Fixes high variance*
- Fixes high bias*

Diagnosing our datasets

- Learning curves can be also used to diagnose the quality of our training/validation sets

Unrepresentative Training Set

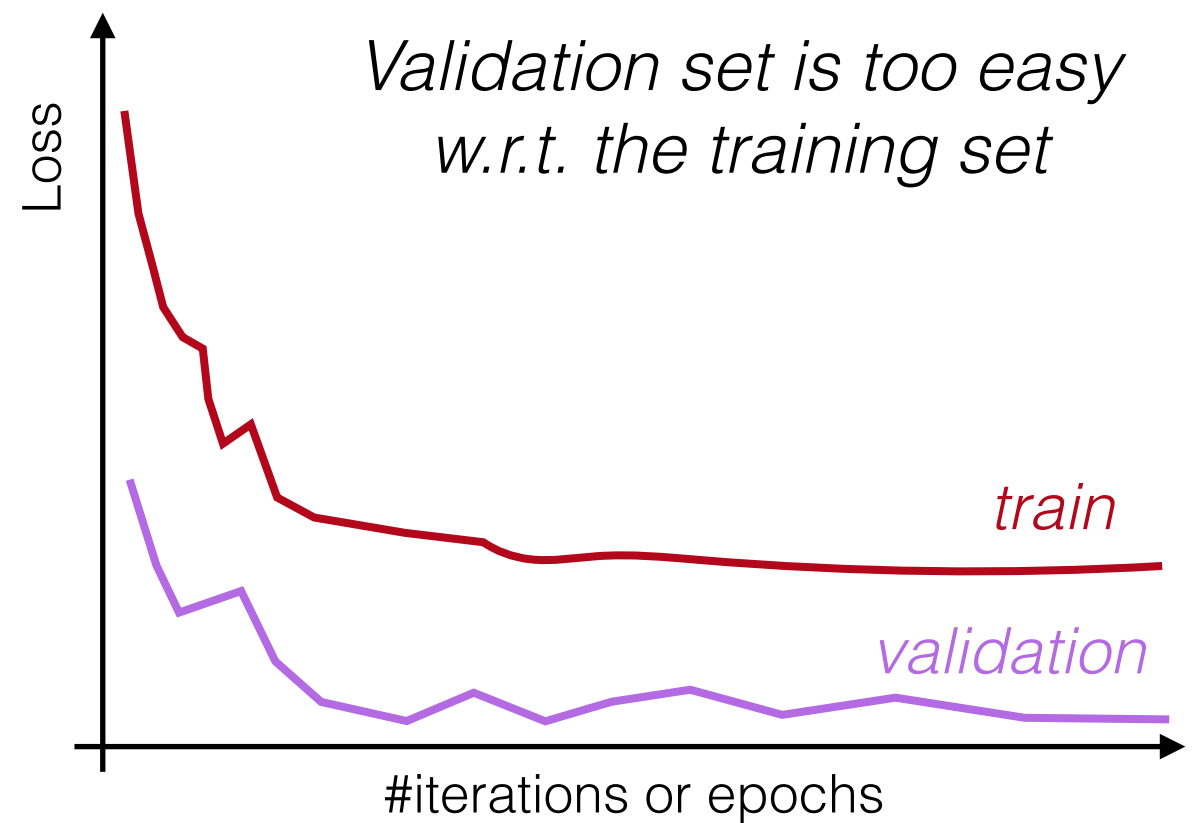
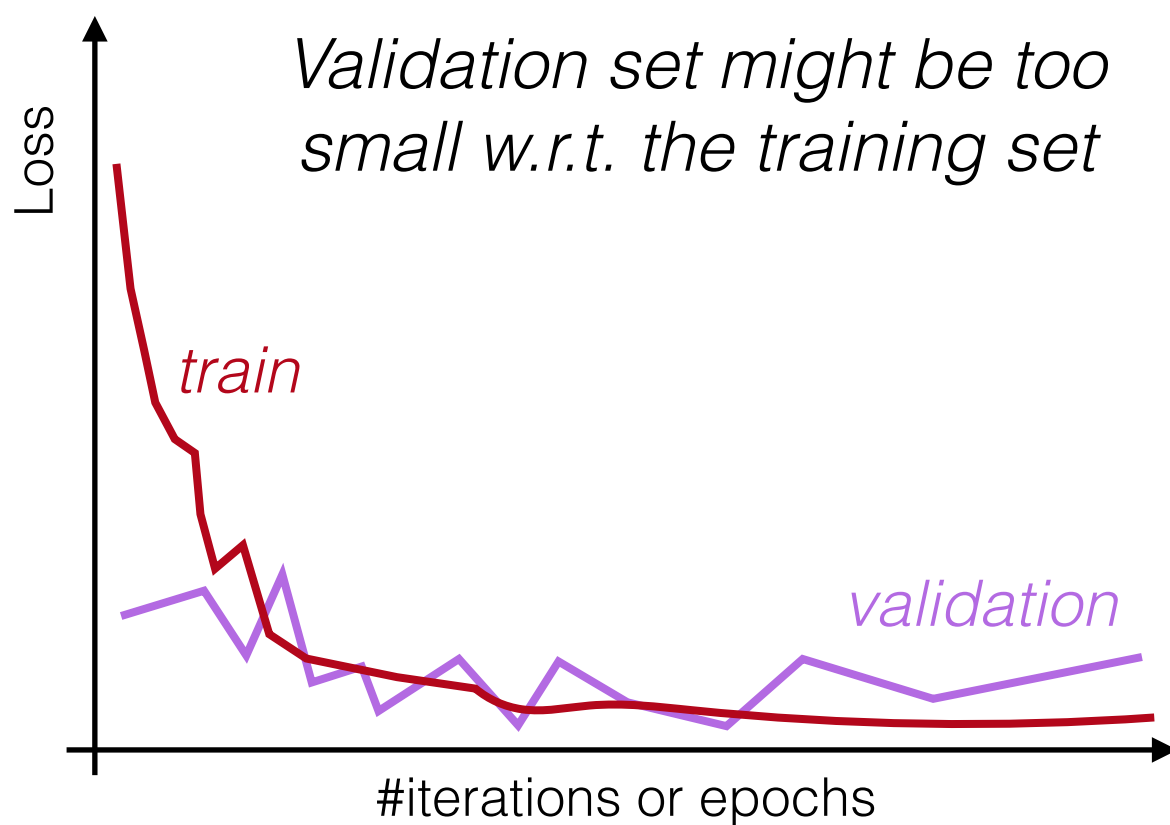


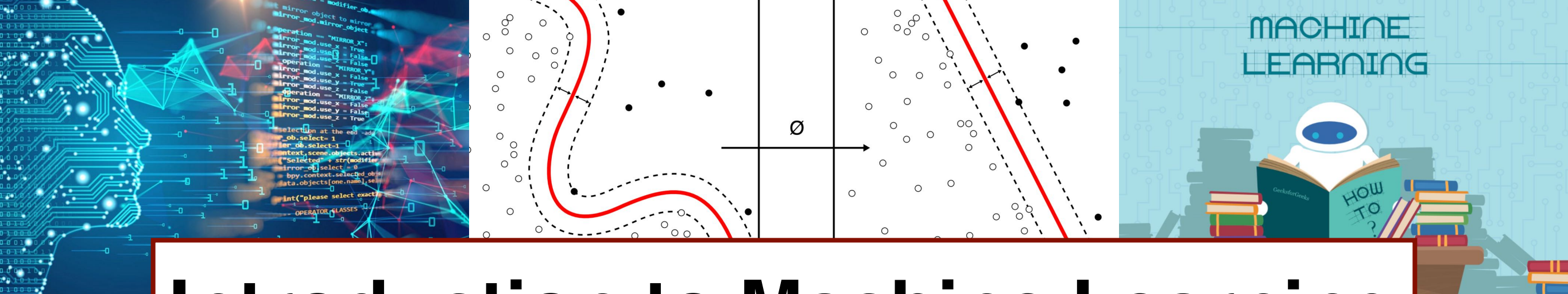
- ▶ The training set does not provide sufficient information to learn the problem
- ▶ It may occur if the training set has too few examples as compared to the validation set

Diagnosing our datasets

- Learning curves can be also used to diagnose the quality of our training/validation sets

Unrepresentative Validation Set





Introduction to Machine Learning

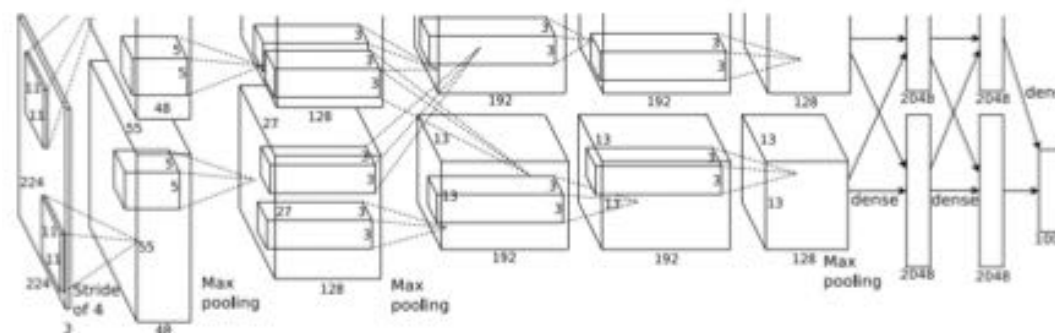
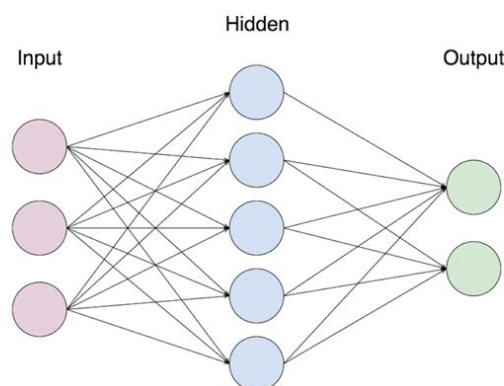
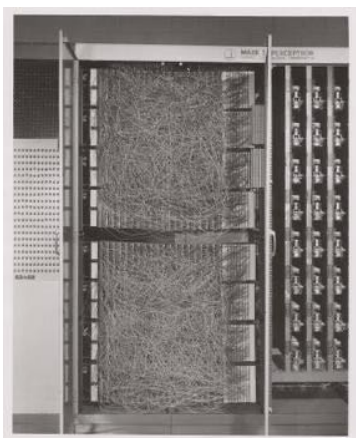
SCP8084699 - LT Informatica



Artificial Neural Networks I
Prof. Lamberto Ballan

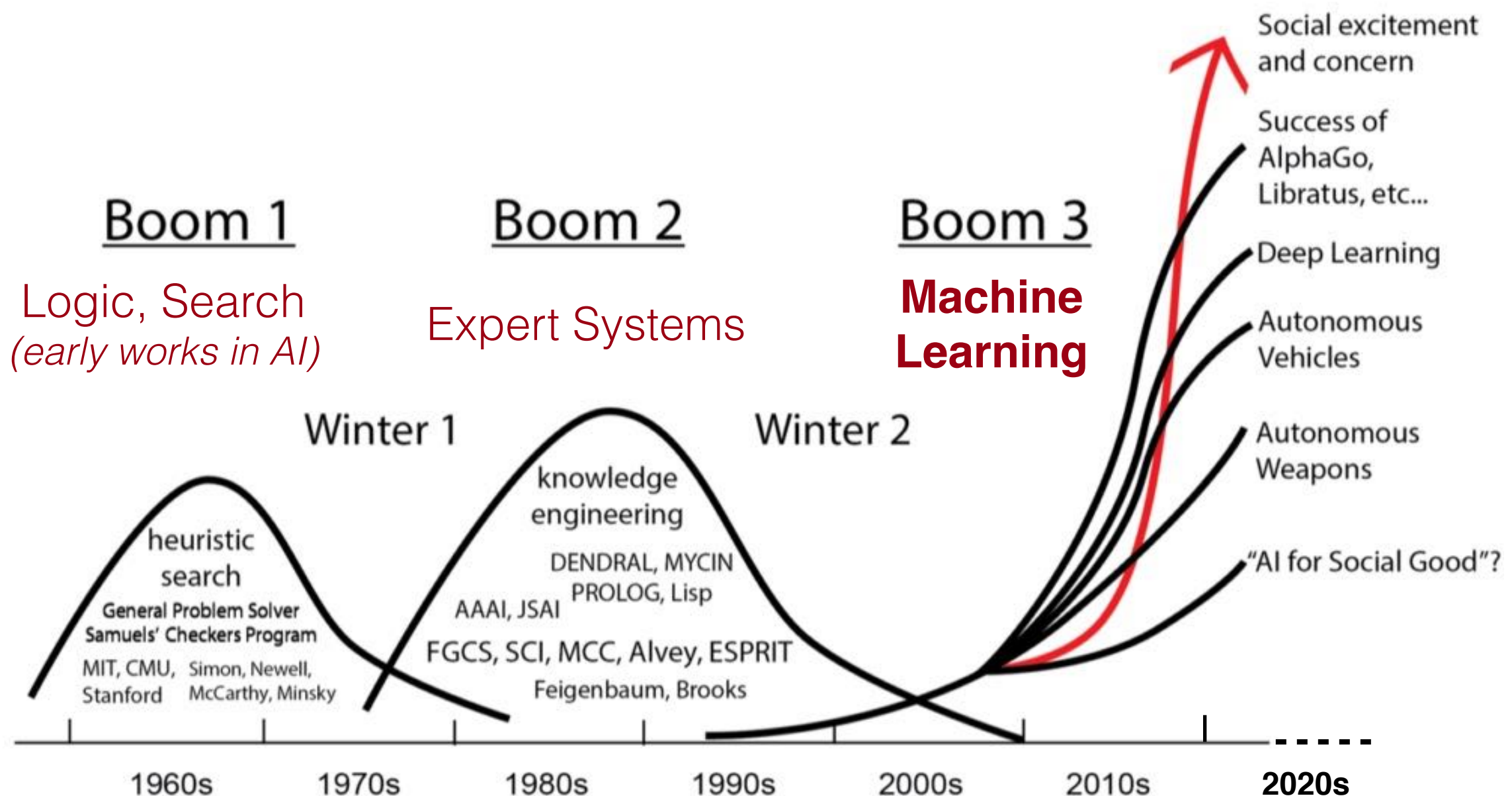
Neural Networks: origins

- Algorithms that are modelled (loosely) after the brain
 - Neural networks have been introduced in late 1950s
 - They were widely used in 1980s and early 1990s; their popularity largely diminished in late 1990s
 - Recent “resurgence”: late 2000s/early 2010s (deep learning revolution)
 - Artificial neural networks are not nearly as complex or intricate as the actual brain structure



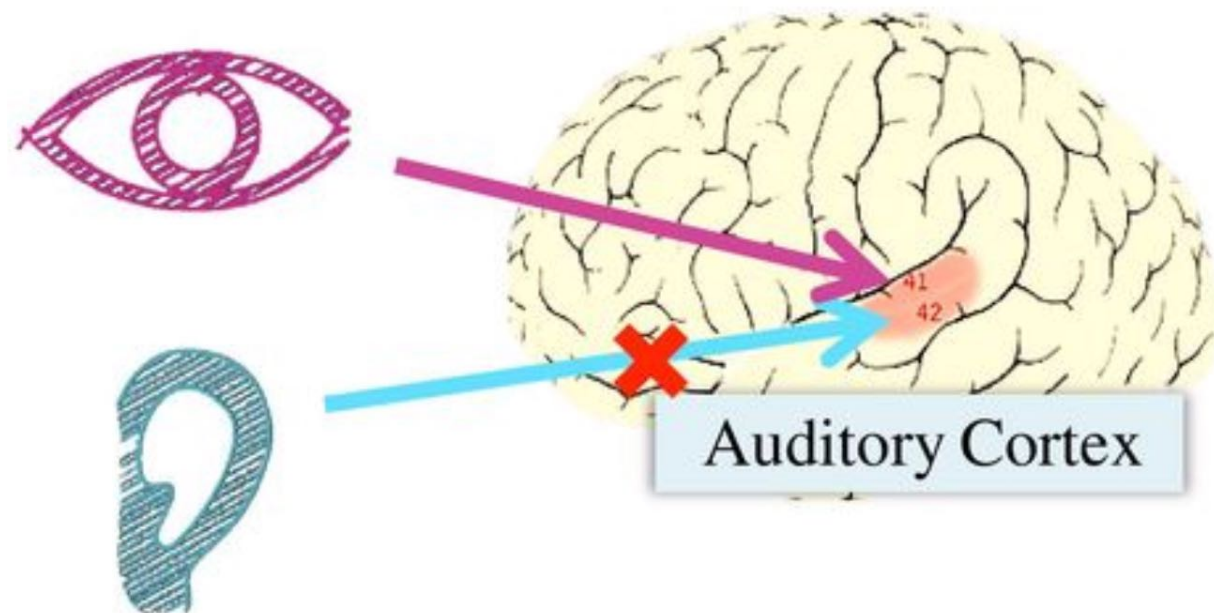
Neural Networks: origins

- Neural Networks and AI winters



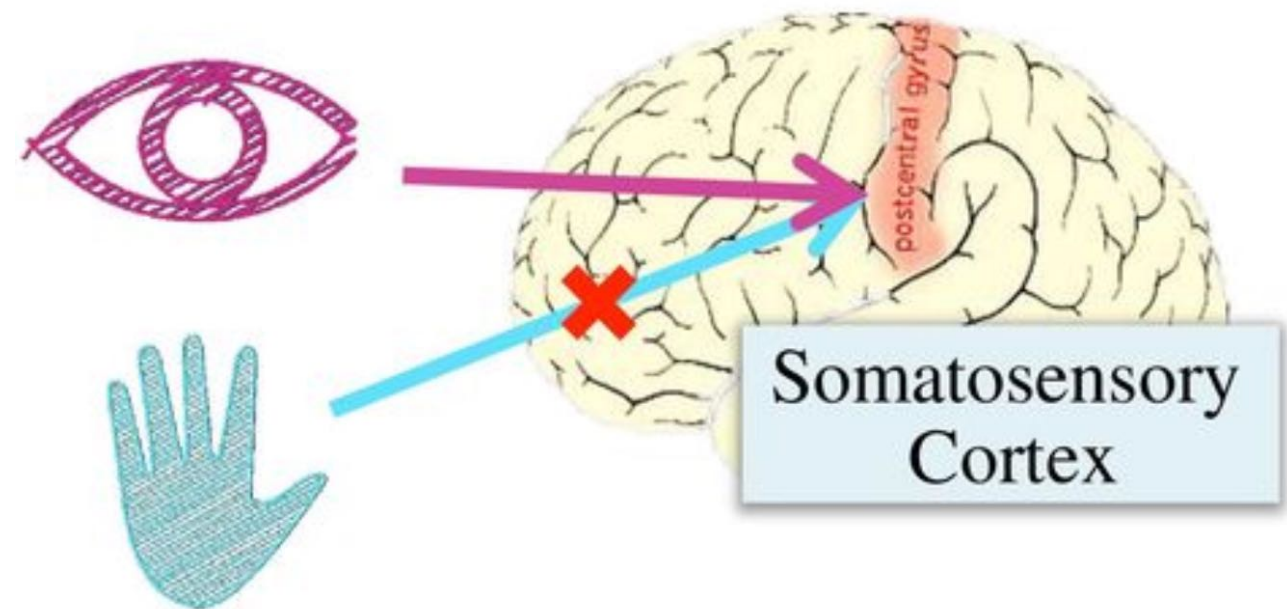
Neural Networks: origins

- The “one learning algorithm” hypothesis



Auditory cortex learns to see

[Roe et al., 1992]

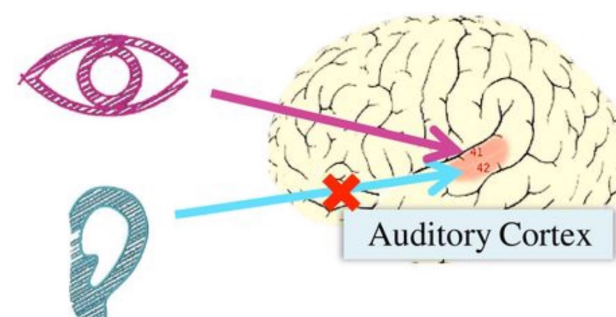


Somatosensory cortex
learns to see

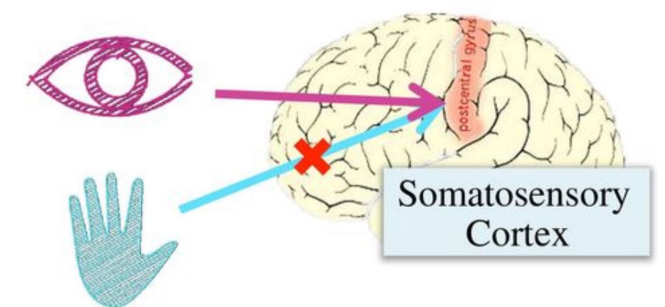
[Metin and Frost, 1989]

Neural Networks: origins

- The “one learning algorithm” hypothesis
 - ▶ Because of these experiments there is this sense that if the same part of the brain can process sight or sound or touch, then (maybe) there is *one learning algorithm*
 - ▶ Therefore, instead of designing hundreds of different algorithms, we need to approximate this one algorithm



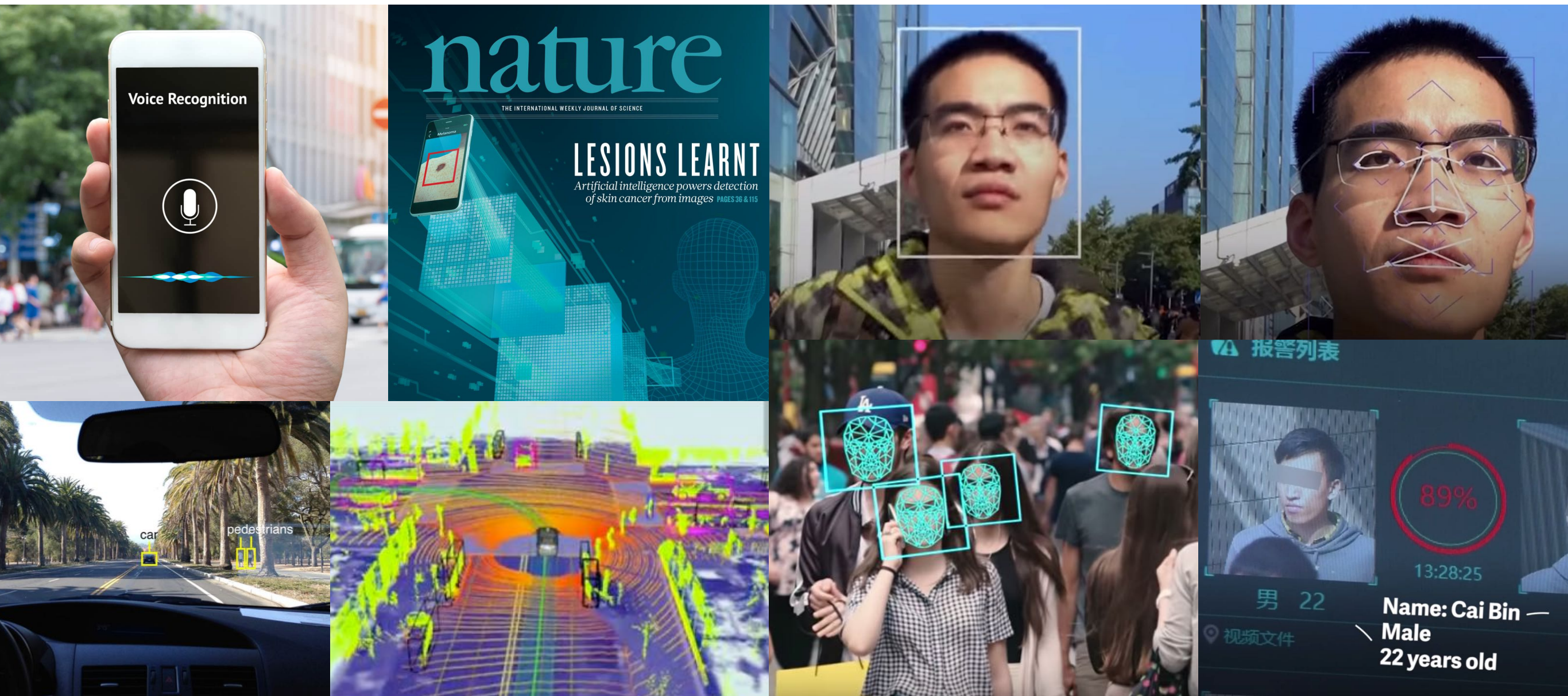
[Roe et al., 1992]



[Metin and Frost, 1989]

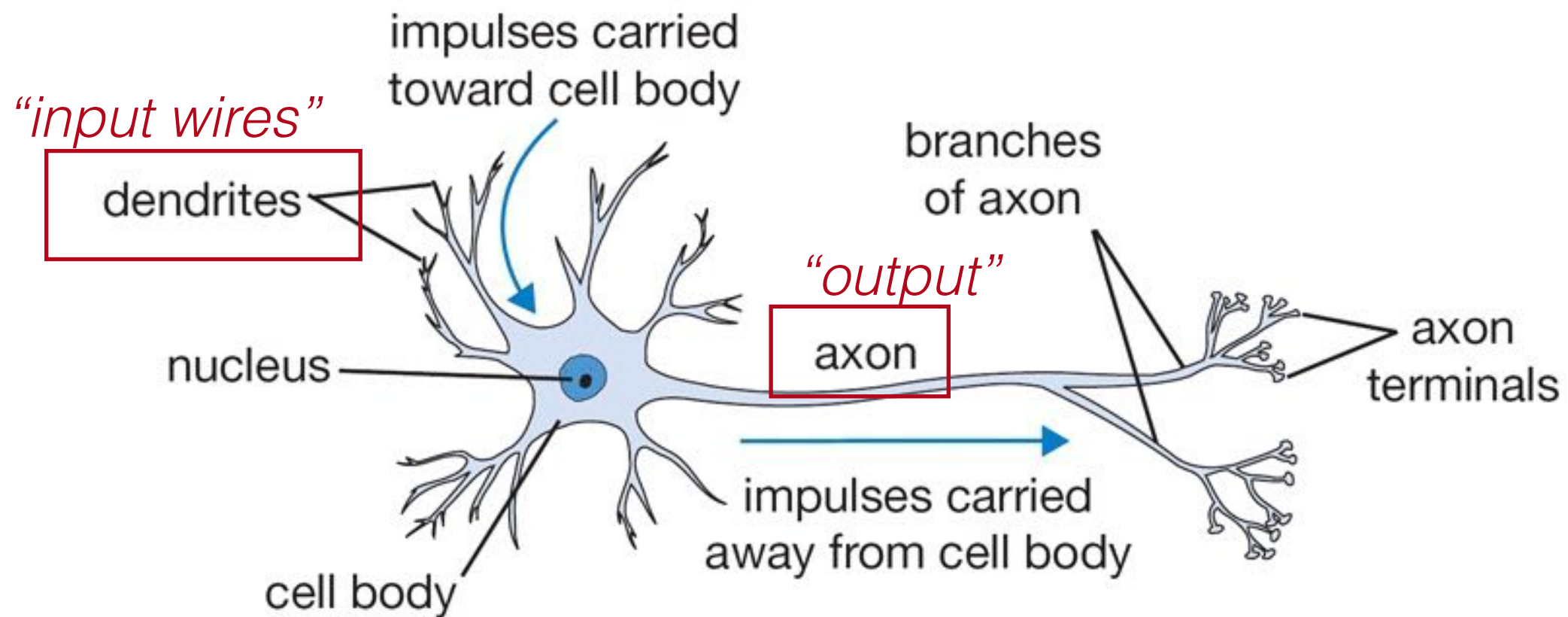
Neural Networks: applications

- Neural networks are everywhere...



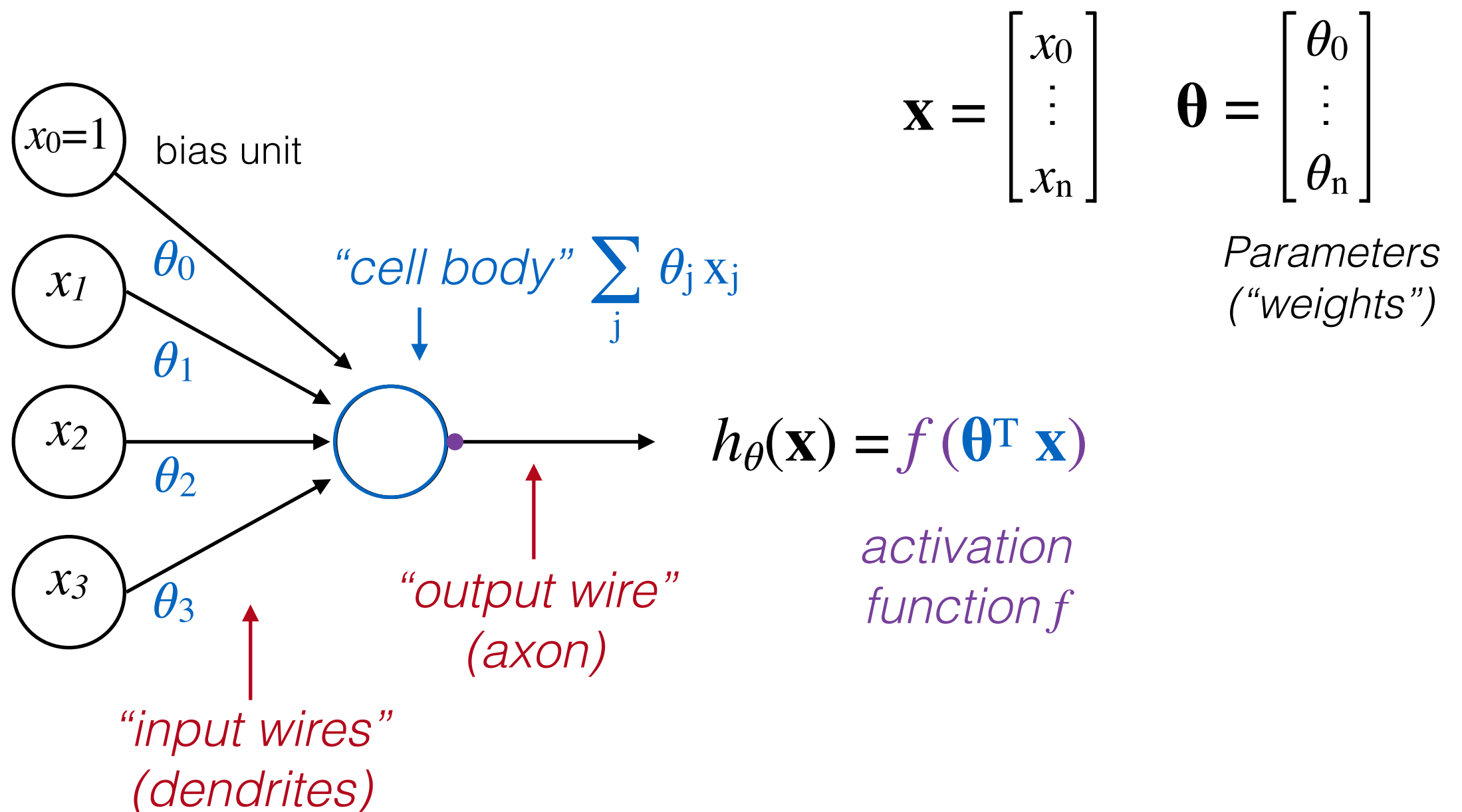
Neural Networks

- Let's start by looking at what a single neuron in the brain looks like



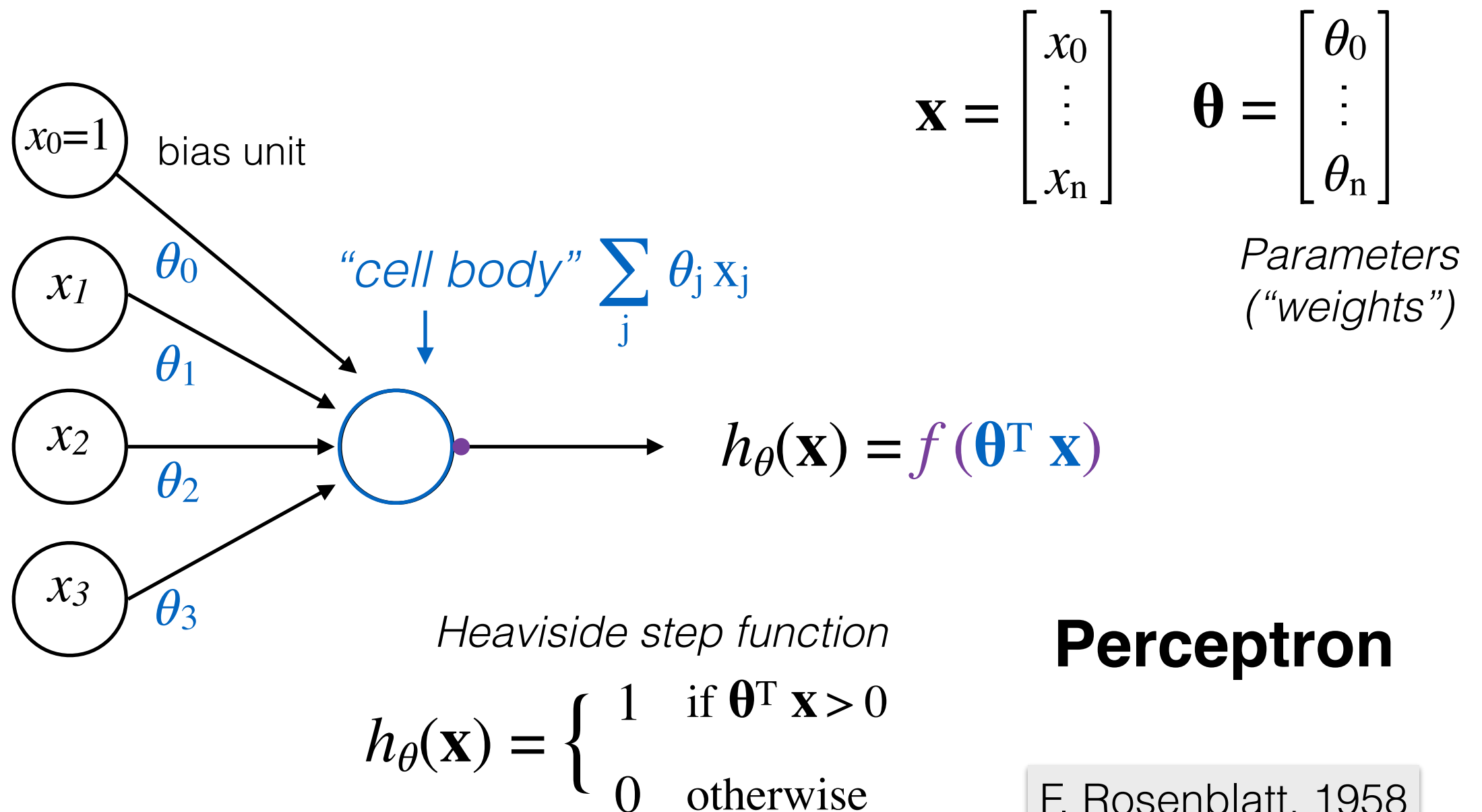
Artificial Neural Networks

- We define a simple model for an artificial neuron (computational unit)



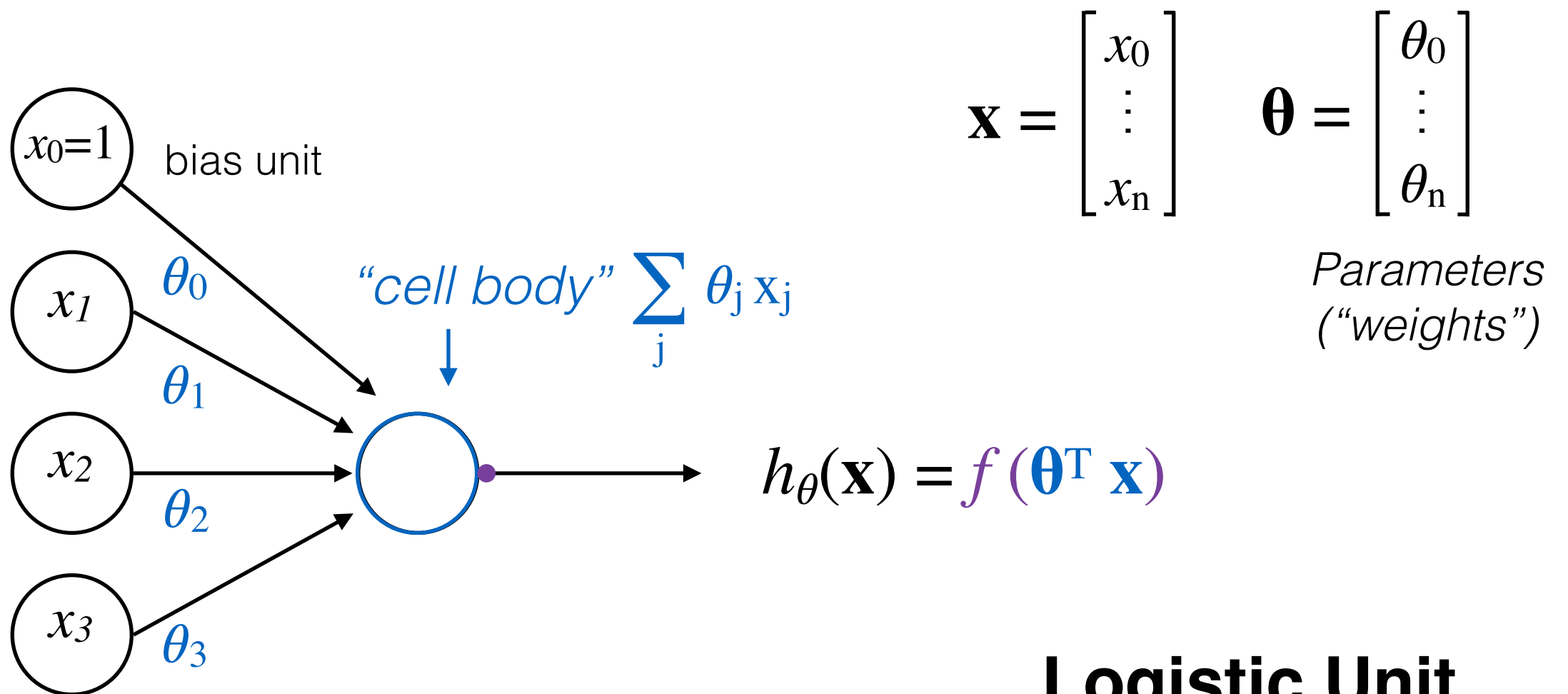
Artificial Neural Networks

- We define a simple model for an artificial neuron (computational unit)



Artificial Neural Networks

- We define a simple model for an artificial neuron (computational unit)



Logistic Unit

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

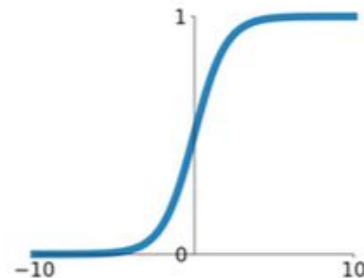
"activation function"
"non-linearity"

Neural Networks: activation functions

- Nowadays there are several activation functions you may encounter:

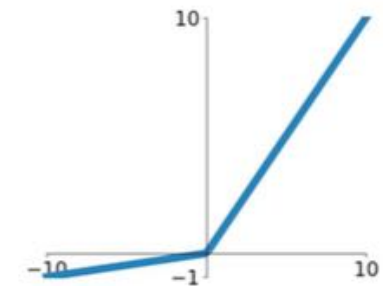
Sigmoid

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



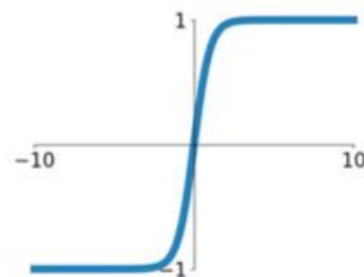
Leaky ReLU

$$\max(0.1x, x)$$



tanh

$$\tanh(x)$$

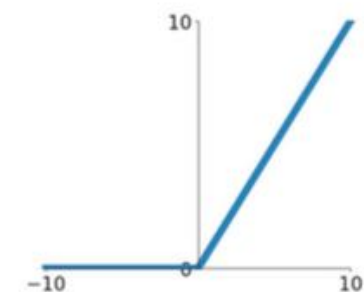


Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

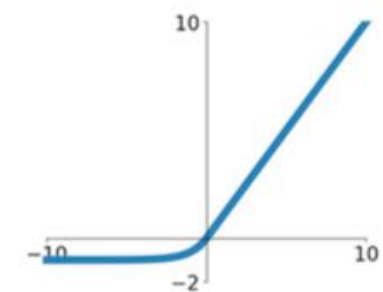
ReLU

$$\max(0, x)$$



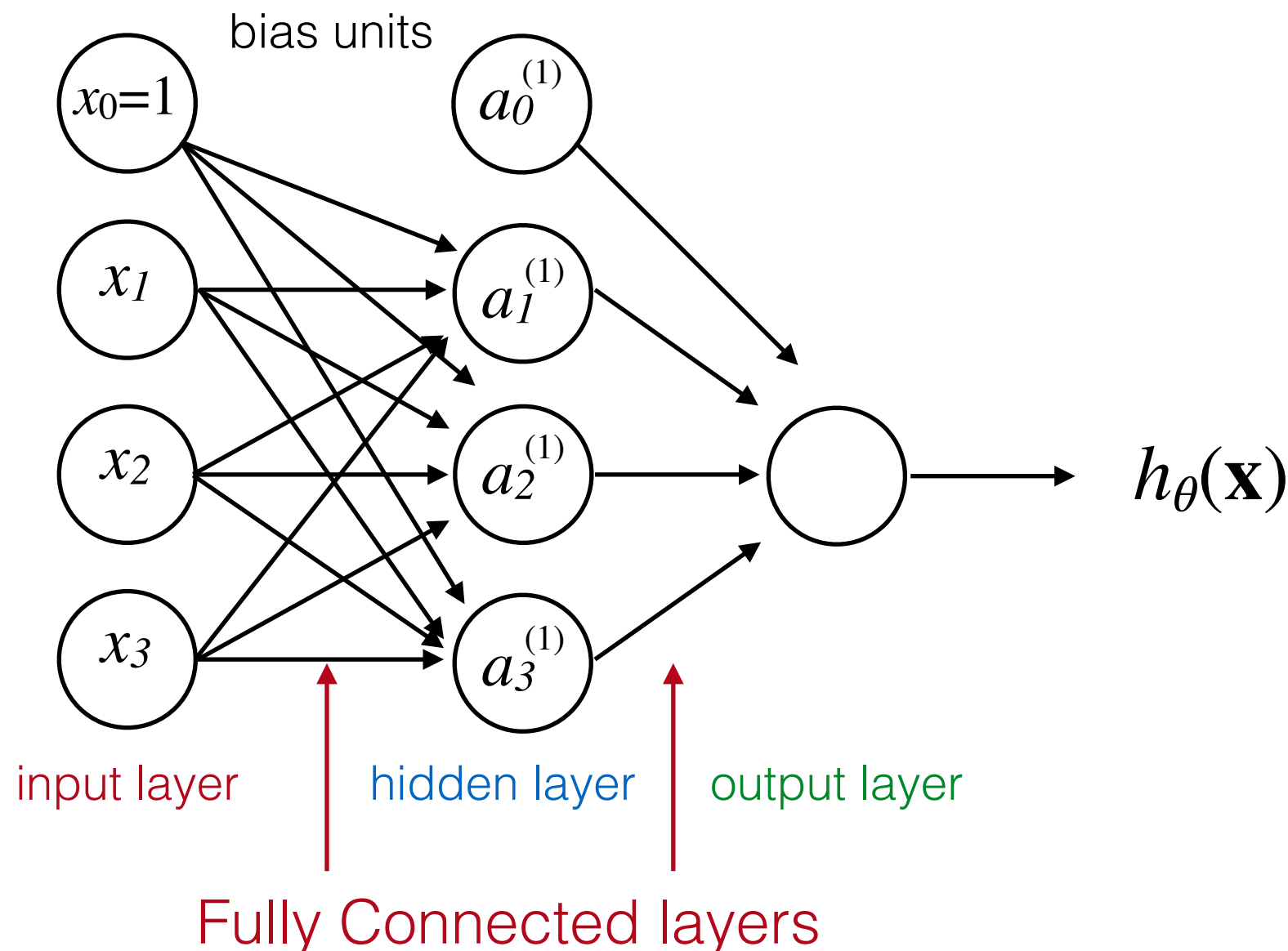
ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



Neural Networks: definition

- A (artificial) neural network is just a group of this different neurons strong together

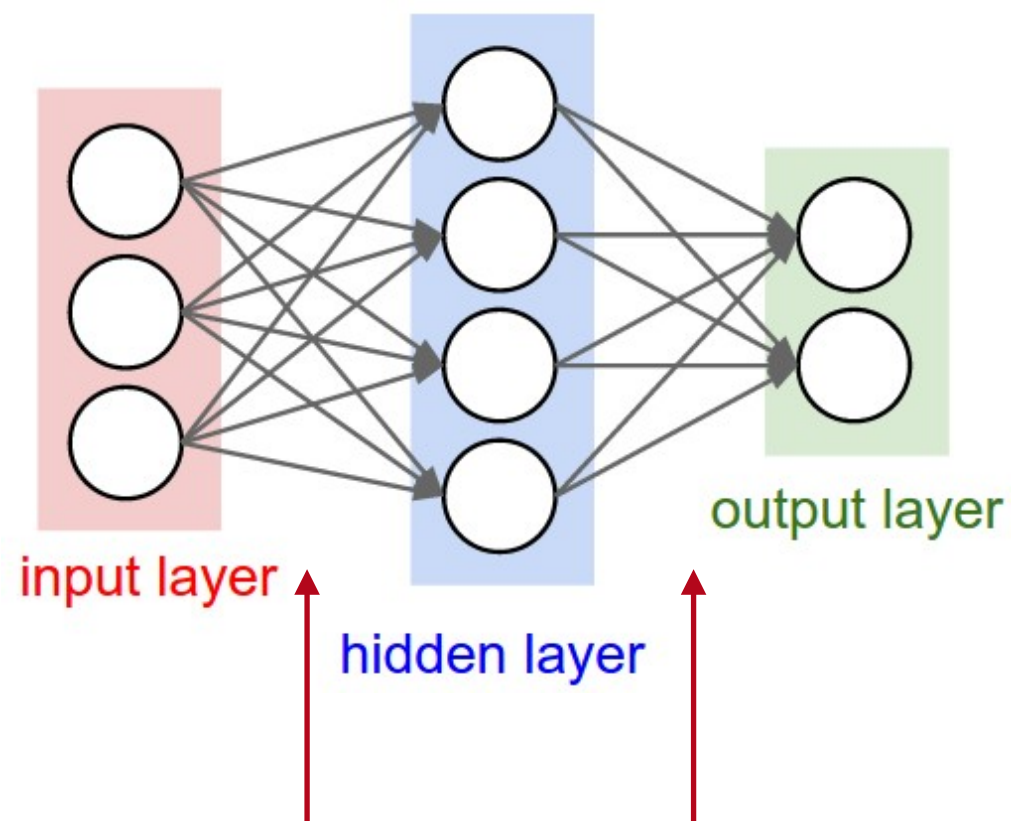


Neural Networks: architectures

- Here the term *architecture* refers to how the different neurons are connected to each other

2-layer neural network

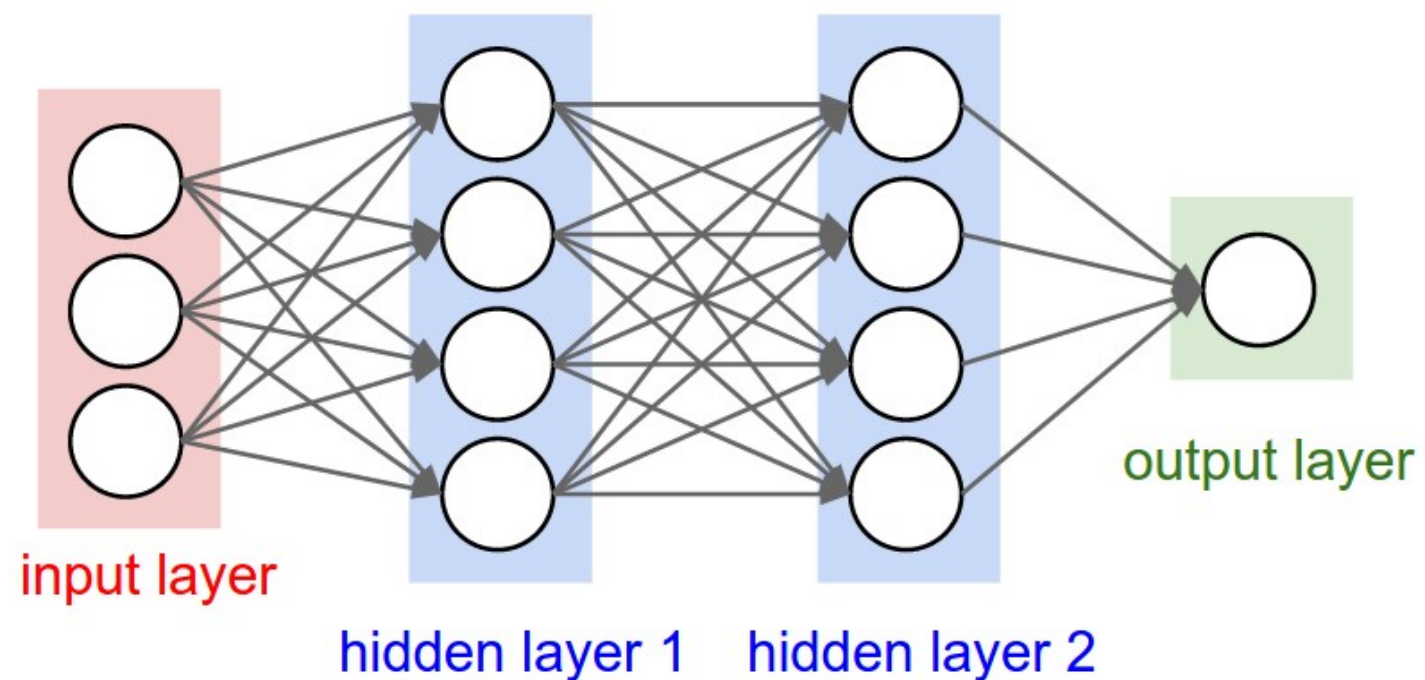
or 1-hidden-layer neural net



Fully Connected layers

3-layer neural network

or 2-hidden-layer neural net

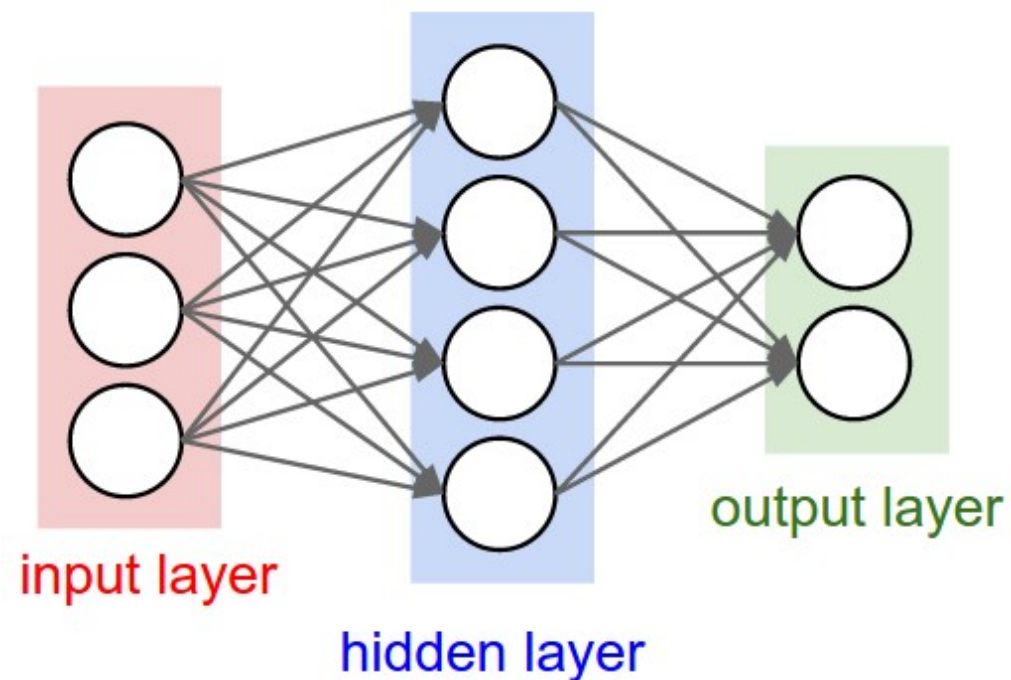


Note: when we say *N-layer neural network*, we do not count the input layer

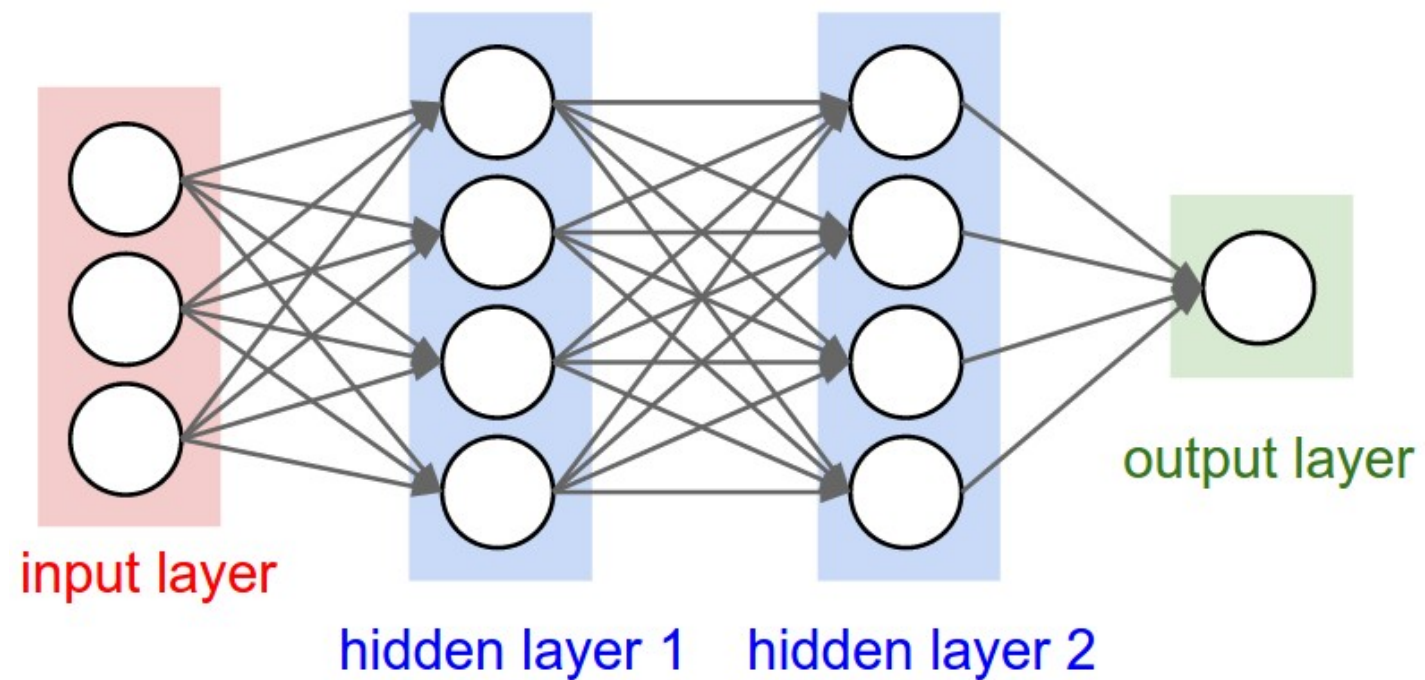
Neural Networks: architectures

- **Sizing neural networks:** the two metrics that are commonly used to measure the size of a NeuralNet are the #neurons or (more commonly) the #parameters

(a) **2-layer neural network**



(b) **3-layer neural network**

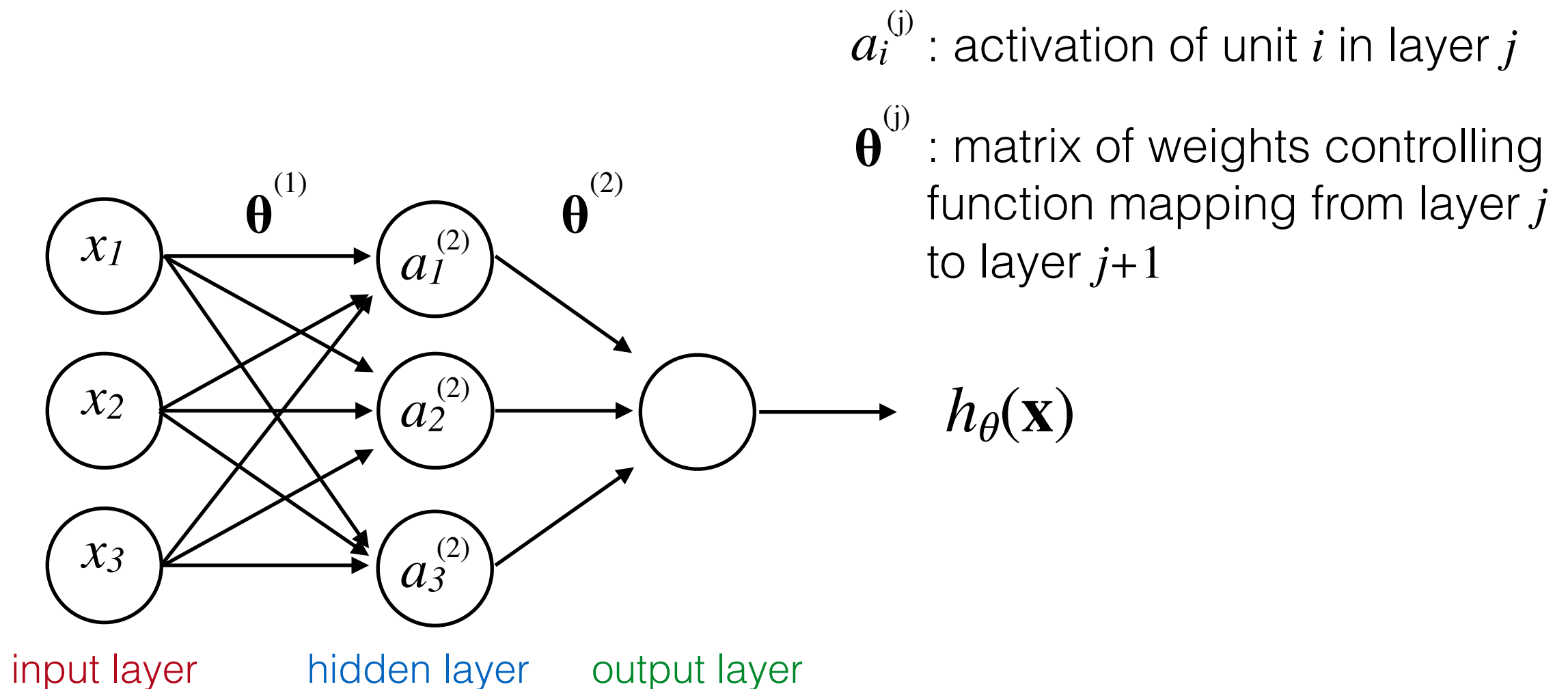


(b) has $4+4+1$ neurons; how many (learnable) parameters?

(b) #parameters: $3 \times 4 + 4 \times 4 + 4 \times 1 = 32 + 9$ (*biases*) = 41

Neural Networks: model representation

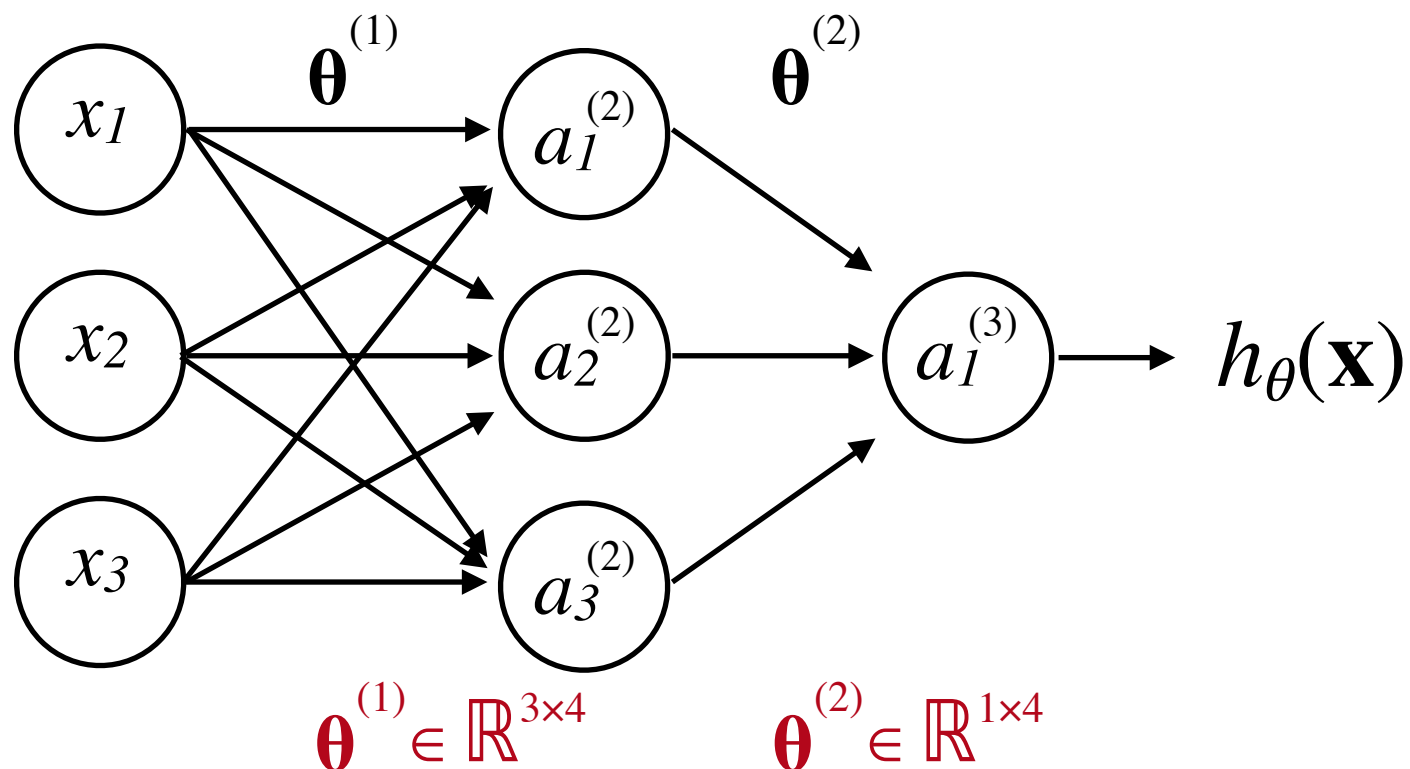
- A (artificial) neural network is just a group of this different neurons strong together



Neural Networks: feed-forward computation

- Input units are set by some exterior function (it can be generated by sensors) which causes their output links to be activated at the specified level
- Working forward through the network, these outputs are going to be the input for the next layer
 - Each output is just the weighted sum of the activation on the links feeding into a node
 - The activation function transforms this linear combination: typically this is a non linear function (such a sigmoid)
 - This function corresponds to the “threshold” of that node

Neural Networks: feed-forward computation



$a_i^{(j)}$: activation of unit i in layer j
 $\theta^{(j)}$: matrix of weights controlling function mapping from layer j to layer $j+1$

If network has u_j units in layer j and u_{j+1} units in layer $j+1$, then $\theta^{(j)}$ will be of dimension $u_{j+1} \times (u_j + 1)$

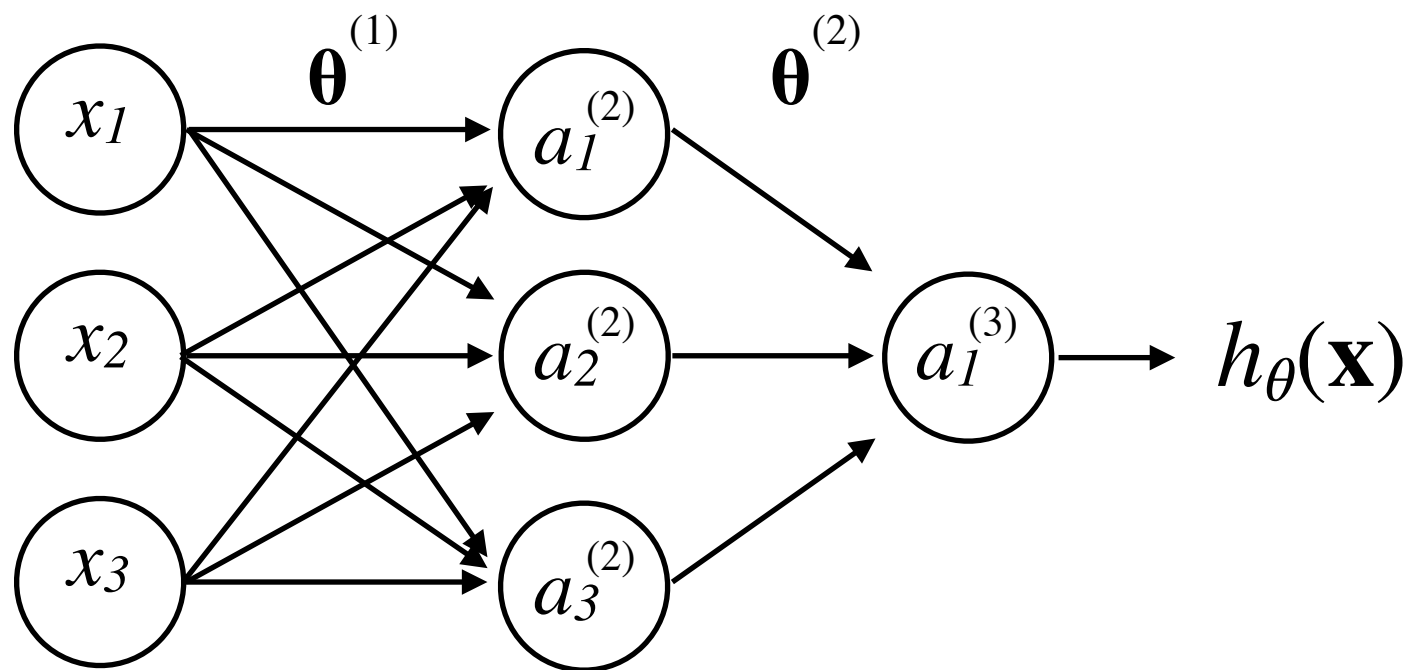
$$a_1^{(2)} = f(\theta_{10}^{(1)} x_0 + \theta_{11}^{(1)} x_1 + \theta_{12}^{(1)} x_2 + \theta_{13}^{(1)} x_3)$$

$$a_2^{(2)} = f(\theta_{20}^{(1)} x_0 + \theta_{21}^{(1)} x_1 + \theta_{22}^{(1)} x_2 + \theta_{23}^{(1)} x_3)$$

$$a_3^{(2)} = f(\theta_{30}^{(1)} x_0 + \theta_{31}^{(1)} x_1 + \theta_{32}^{(1)} x_2 + \theta_{33}^{(1)} x_3)$$

$$h_{\theta}(\mathbf{x}) = a_1^{(3)} = f(\theta_{10}^{(2)} a_0^{(2)} + \theta_{11}^{(2)} a_1^{(2)} + \theta_{12}^{(2)} a_2^{(2)} + \theta_{13}^{(2)} a_3^{(2)})$$

Neural Networks: feed-forward computation



(vectorization)

$a_i^{(j)}$: activation of unit i in layer j

$\theta^{(j)}$: matrix of weights controlling function mapping from layer j to layer $j+1$

$z_i^{(j)}$: linear combination of the input nodes

$$a_1^{(2)} = f(\overbrace{\theta_{10}^{(1)} x_0 + \theta_{11}^{(1)} x_1 + \theta_{12}^{(1)} x_2 + \theta_{13}^{(1)} x_3}^{z_1^{(2)}}) = f(z_1^{(2)})$$

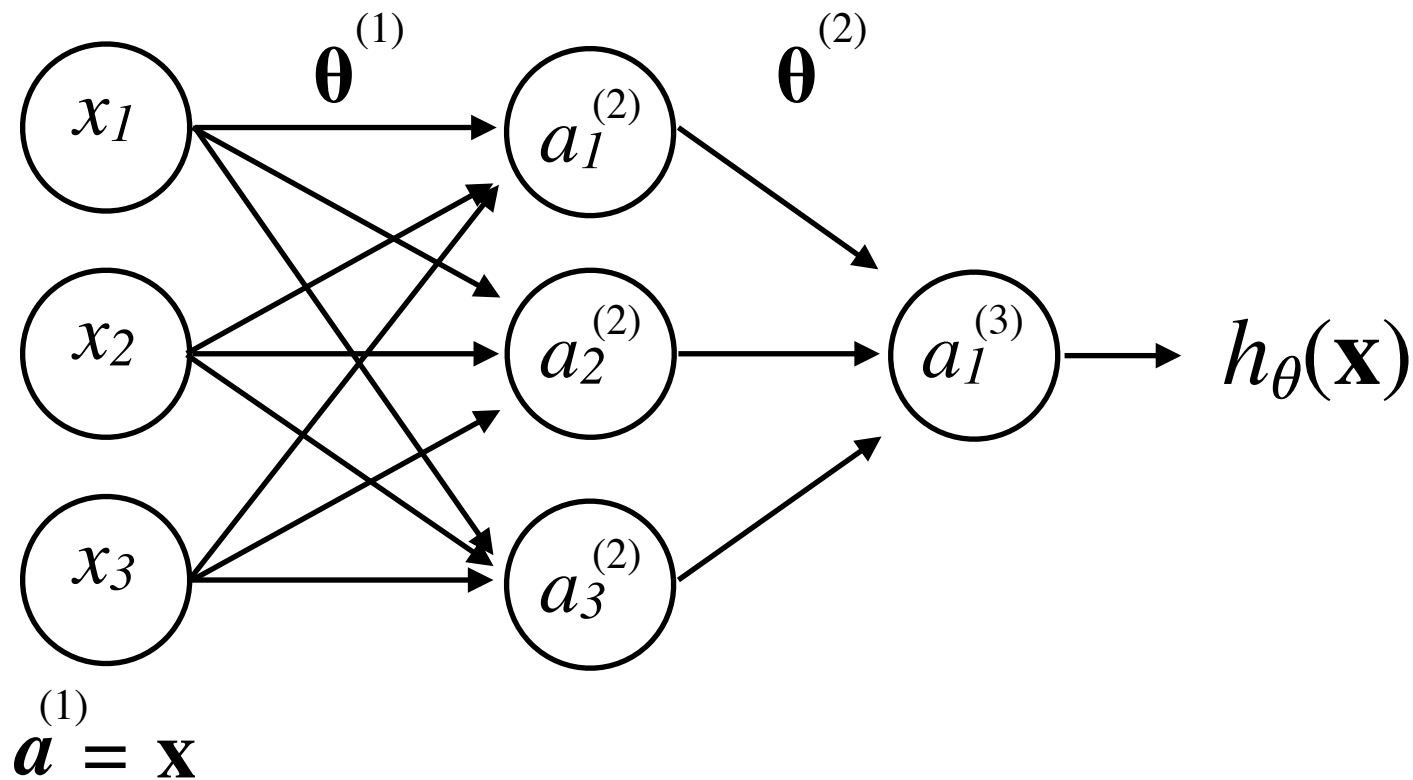
$$a_2^{(2)} = f(\theta_{20}^{(1)} x_0 + \theta_{21}^{(1)} x_1 + \theta_{22}^{(1)} x_2 + \theta_{23}^{(1)} x_3) = f(z_2^{(2)})$$

$$a_3^{(2)} = f(\theta_{30}^{(1)} x_0 + \theta_{31}^{(1)} x_1 + \theta_{32}^{(1)} x_2 + \theta_{33}^{(1)} x_3) = f(z_3^{(2)})$$

$$h_{\theta}(\mathbf{x}) = a_1^{(3)} = f(\theta_{10}^{(2)} a_0^{(2)} + \theta_{11}^{(2)} a_1^{(2)} + \theta_{12}^{(2)} a_2^{(2)} + \theta_{13}^{(2)} a_3^{(2)}) = f(z_1^{(3)})$$

Neural Networks: feed-forward computation

(vectorization)



$$\mathbf{x} = \begin{bmatrix} x_0 \\ \vdots \\ x_3 \end{bmatrix} \quad \mathbf{z}^{(2)} = \begin{bmatrix} z_1^{(2)} \\ \vdots \\ z_3^{(2)} \end{bmatrix}$$

$$a_1^{(2)} = f(\theta_{10}^{(1)} x_0 + \theta_{11}^{(1)} x_1 + \theta_{12}^{(1)} x_2 + \theta_{13}^{(1)} x_3) = f(z_1^{(2)})$$

$$a_2^{(2)} = f(\theta_{20}^{(1)} x_0 + \theta_{21}^{(1)} x_1 + \theta_{22}^{(1)} x_2 + \theta_{23}^{(1)} x_3) = f(z_2^{(2)})$$

$$a_3^{(2)} = f(\theta_{30}^{(1)} x_0 + \theta_{31}^{(1)} x_1 + \theta_{32}^{(1)} x_2 + \theta_{33}^{(1)} x_3) = f(z_3^{(2)})$$

$$\mathbf{z}^{(2)} = \boldsymbol{\theta}^{(1)} \mathbf{a}^{(1)}$$

$$\mathbf{a}^{(2)} = f(\mathbf{z}^{(2)})$$

Add bias: $a_0^{(2)} = 1$

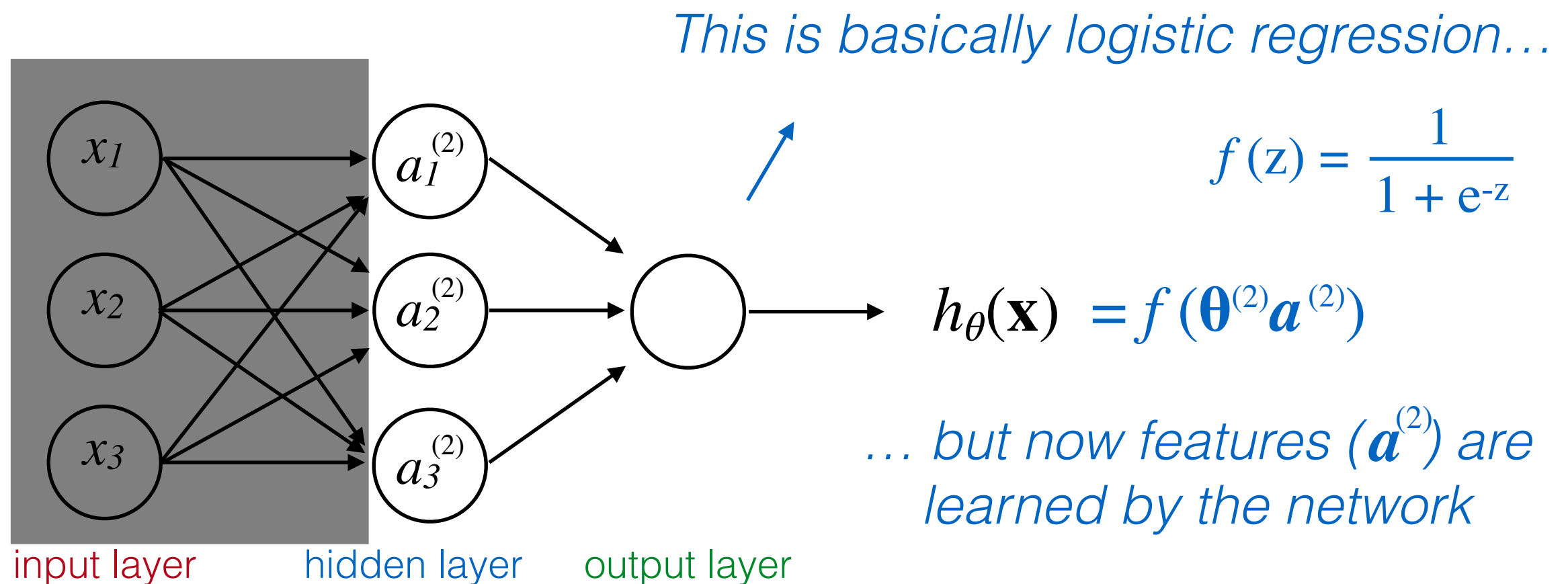
$$\mathbf{z}^{(3)} = \boldsymbol{\theta}^{(2)} \mathbf{a}^{(2)}$$

$$h_{\theta}(\mathbf{x}) = a_1^{(3)} = f(\theta_{10}^{(2)} a_0^{(2)} + \theta_{11}^{(2)} a_1^{(2)} + \theta_{12}^{(2)} a_2^{(2)} + \theta_{13}^{(2)} a_3^{(2)}) = f(z_1^{(3)}) = \mathbf{a}^{(3)} = f(\mathbf{z}^{(3)})$$

“forward propagation”

Neural Networks: feed-forward computation

- This forward propagation view also help us to understand what neural networks might be doing



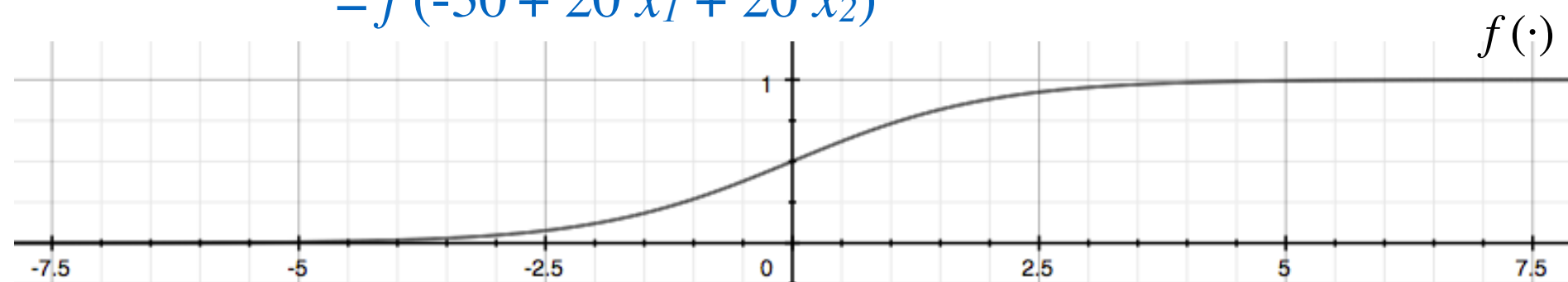
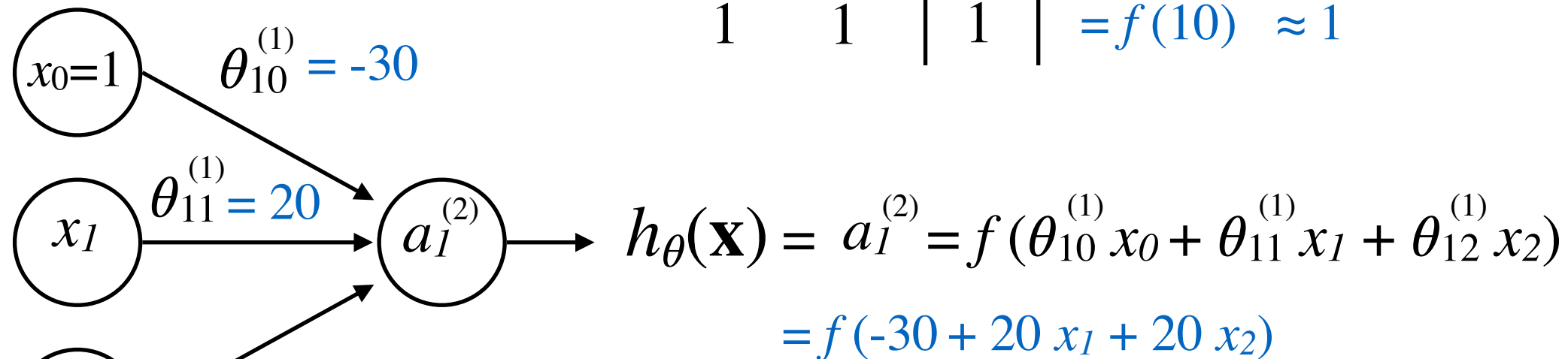
A few examples

- Let's try to compute logical functions

$$x_1, x_2 \in \{0, 1\}$$

$$y = x_1 \text{ AND } x_2$$

x_1	x_2	y	$h_{\theta}(\mathbf{x})$
0	0	0	$=f(-30) \approx 0$
0	1	0	$=f(-10) \approx 0$
1	0	0	$=f(-10) \approx 0$
1	1	1	$=f(10) \approx 1$



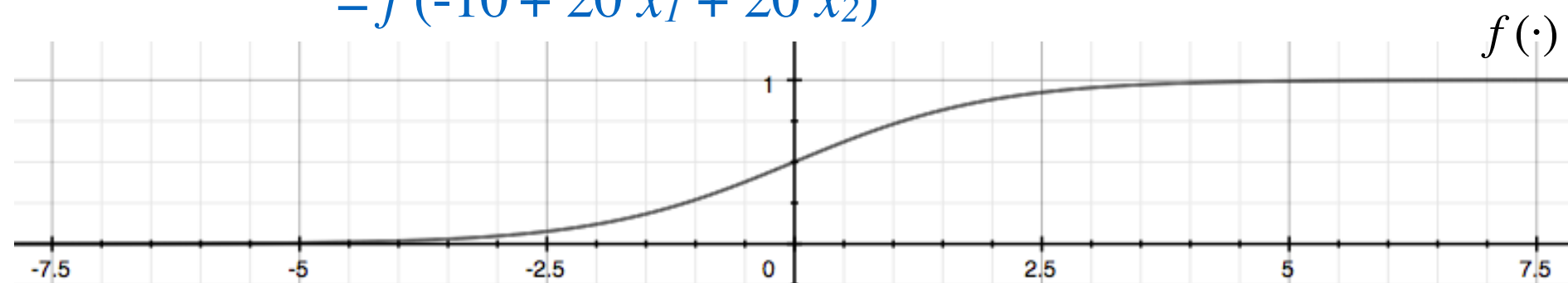
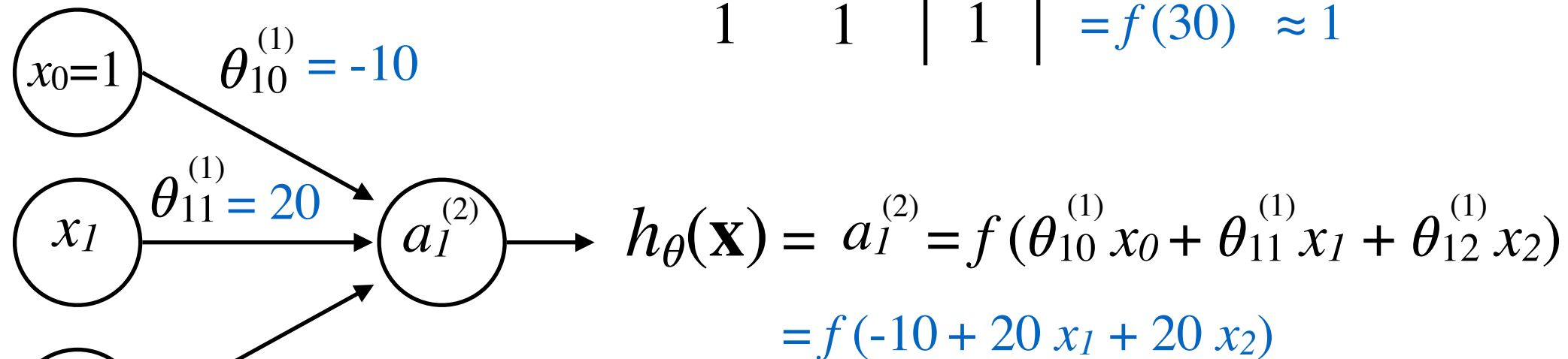
A few examples

- Let's try to compute logical functions

$$x_1, x_2 \in \{0, 1\}$$

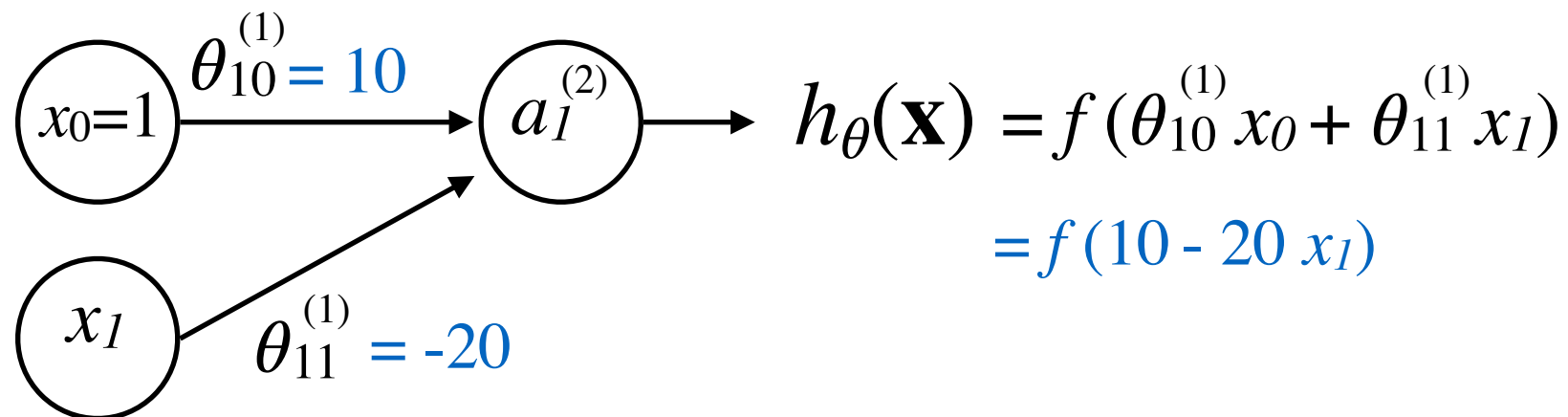
$$y = x_1 \text{ OR } x_2$$

x_1	x_2	y	$h_{\theta}(\mathbf{x})$
0	0	0	$= f(-10) \approx 0$
0	1	1	$= f(10) \approx 1$
1	0	1	$= f(10) \approx 1$
1	1	1	$= f(30) \approx 1$



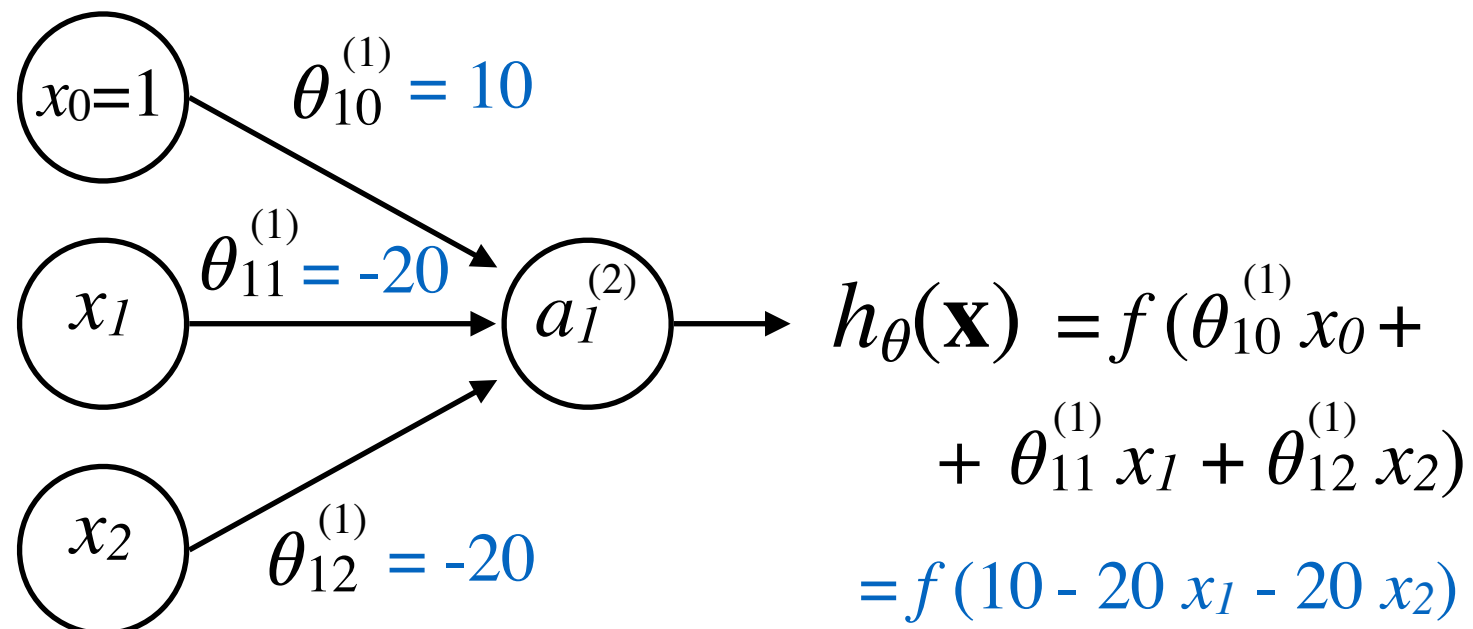
Representing Boolean functions

$$y = (\text{NOT } x_1)$$



x_1	y	$h_{\theta}(\mathbf{x})$
0	1	$= f(10) \approx 1$
1	0	$= f(-10) \approx 0$

$$y = (\text{NOT } x_1) \text{ AND } (\text{NOT } x_2)$$



x_1	x_2	y	$h_{\theta}(\mathbf{x})$
0	0	1	$= f(10) \approx 1$
0	1	0	$= f(-10) \approx 0$
1	0	0	$= f(-10) \approx 0$
1	1	0	$= f(-30) \approx 0$

Representing Boolean functions

- In 1969 Minsky and Seymour showed that it was impossible for perceptrons to learn an XOR function
 - ▶ Often it's (incorrectly) reported that they also conjectured that a similar result would hold for multi-layer perceptrons
 - ▶ Nevertheless, it caused a significant decline in interest and funding of neural network research
 - ▶ It took ten more years until neural networks experienced a resurgence in the 1980s

M.Minsky, S.Papert, "Perceptrons: an introduction to computational geometry", MIT Press, 1969 (*expanded edition published in 1987*)

Contact

- **Office:** Torre Archimede 6CD, room 622
- **Office hours** (ricevimento): Monday 11:00-13:00

✉ lamberto.ballan@unipd.it

🏠 <http://www.lambertoballan.net>

🏠 <http://vimp.math.unipd.it>

@ twitter.com/lambertoballan