

Department of Computer Science
University of Bristol

COMSM0045 – Applied Deep Learning
comsm0045-applied-deep-learning.github.io

2020/21

Lecture 01

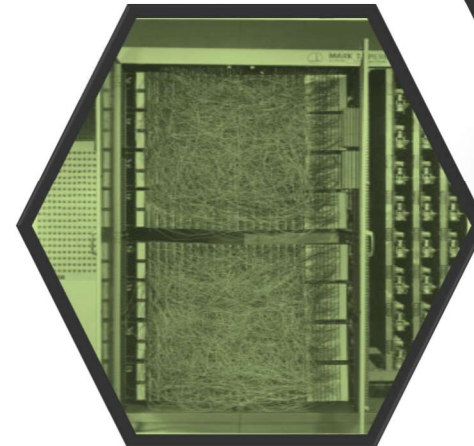
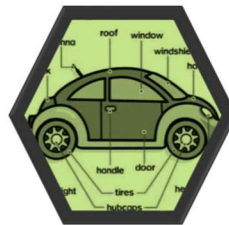
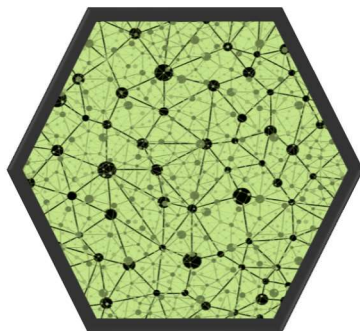
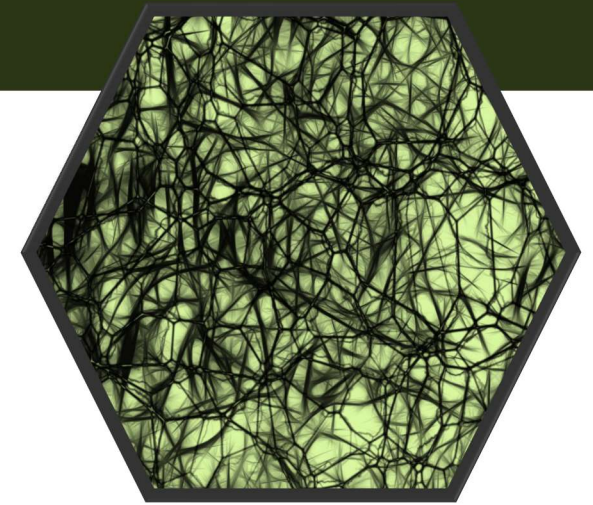
BASICS OF ARTIFICIAL NEURAL NETWORKS

Tilo Burghardt | tilo@cs.bris.ac.uk

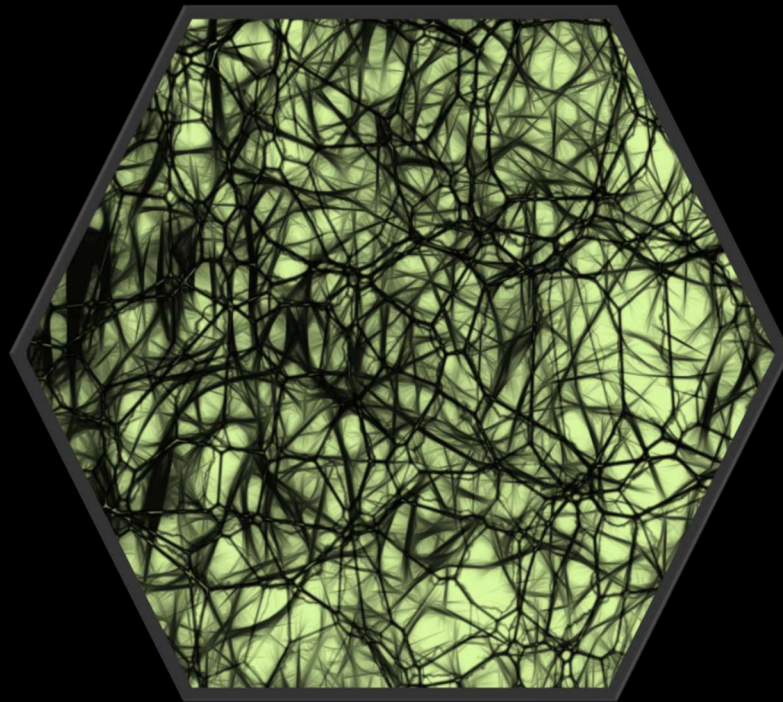
35 Slides

Agenda for Lecture 1

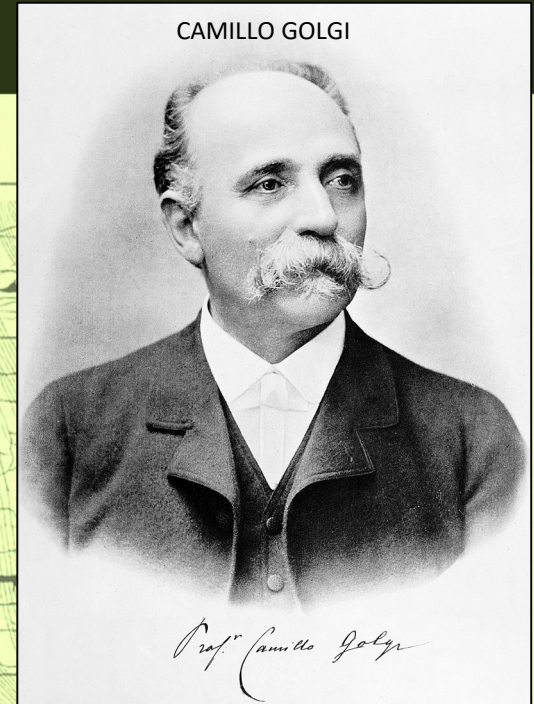
- Neurons and their Structure
- Single & Multi Layer Perceptron
- Basics of Cost Functions
- Gradient Descent and Delta Rule
- Notation and Structure of Deep Feed-Forward Networks



BIOLOGICAL INSPIRATION

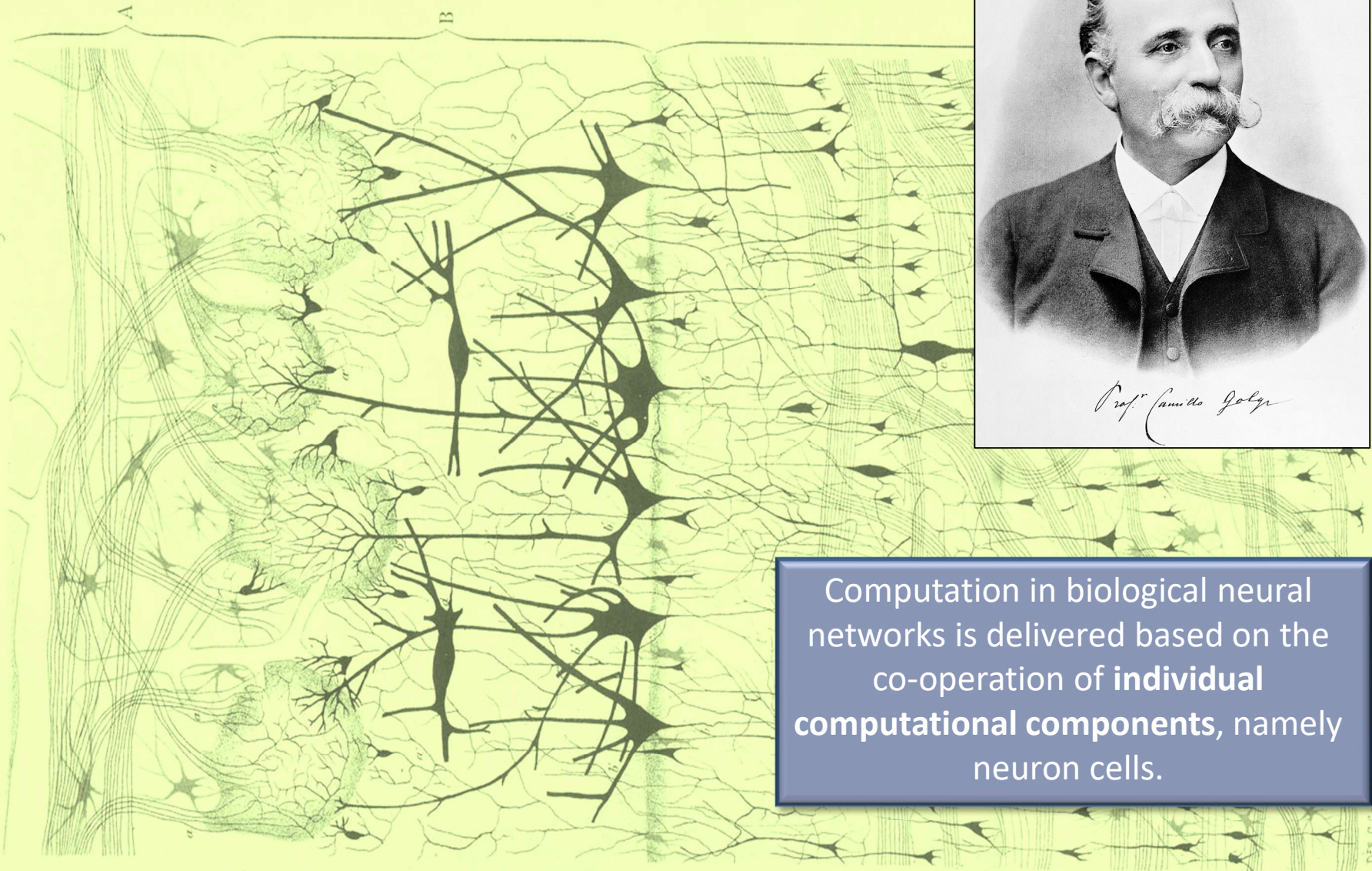


Golgi's first Drawings of Neurons



D^{re} C. Golgi - Bulbi olfattorii

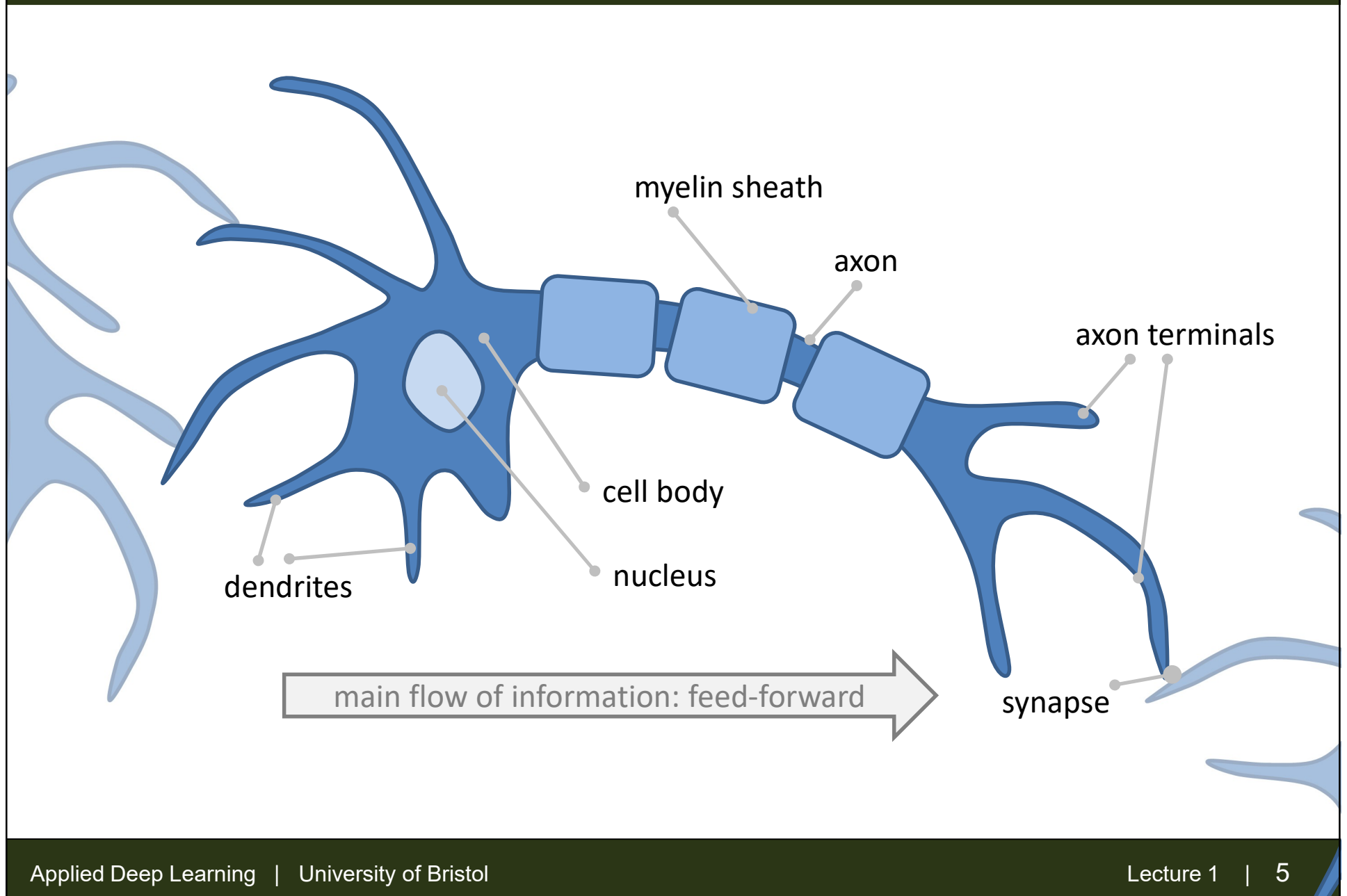
Tav. VII



Computation in biological neural networks is delivered based on the co-operation of **individual computational components**, namely neuron cells.

image source: www.the-scientist.com

Schematic Model of a Neuron



Pavlov and Assistant Conditioning a Dog

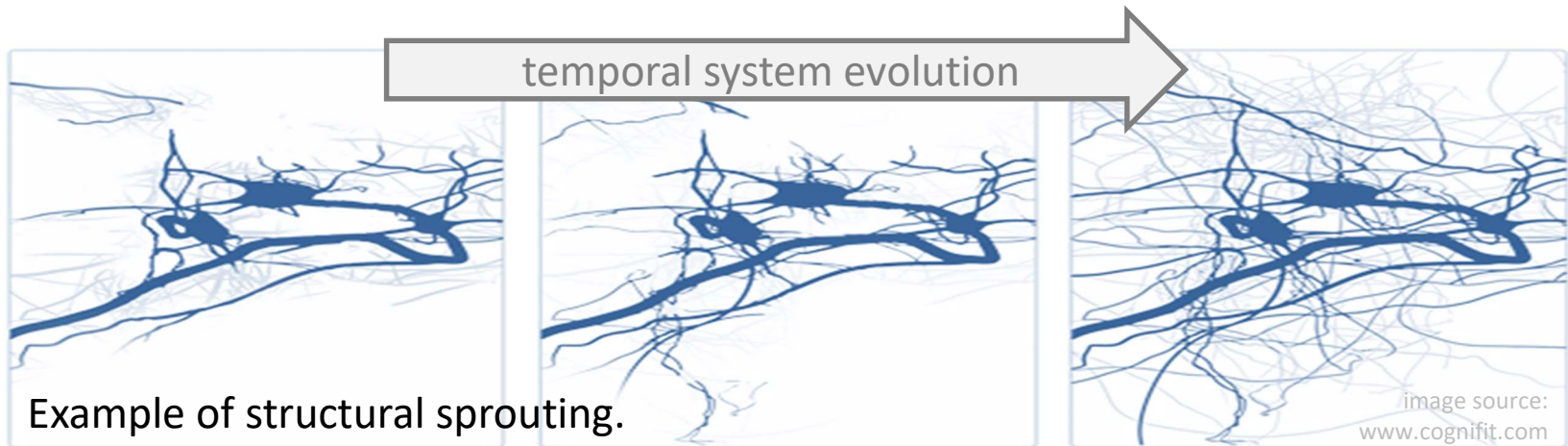


An environment can **condition** the behaviour of biological neural networks leading to the incorporation of new information.

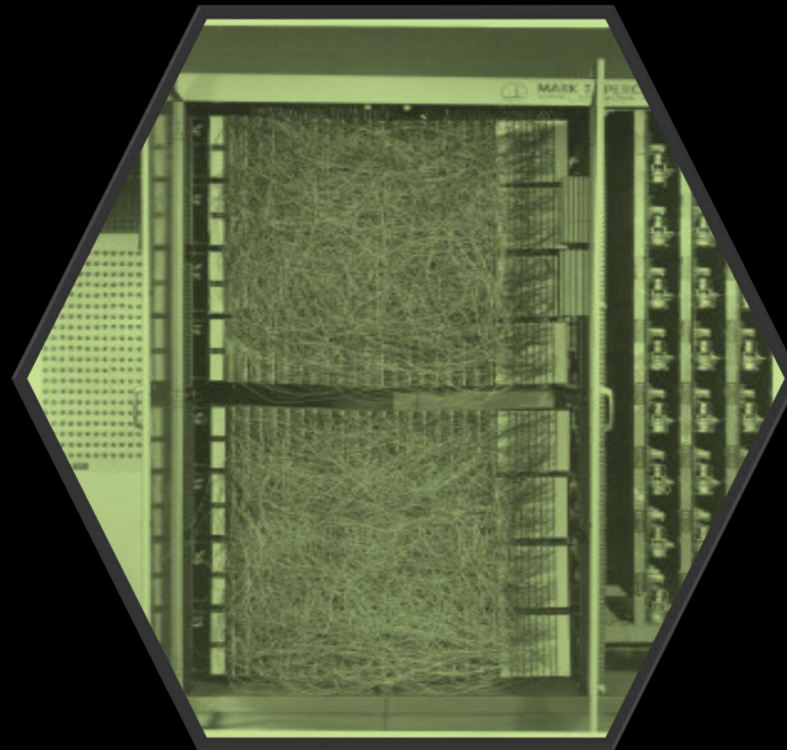
image source: www.psychi.co

Neuro-Plasticity

- plasticity refers to a system's ability to adapt structure and/or behaviour to accommodate new information
- the brain shows various forms of plasticity:
 - natural forms include synaptic plasticity (mainly chemical), structural sprouting (growth), rerouting (functional changes), and neurogenesis (new neurons)



ARTIFICIAL FEED-FORWARD NETWORKS



Rosenblatt's (left) development of the Perceptron (1950s)

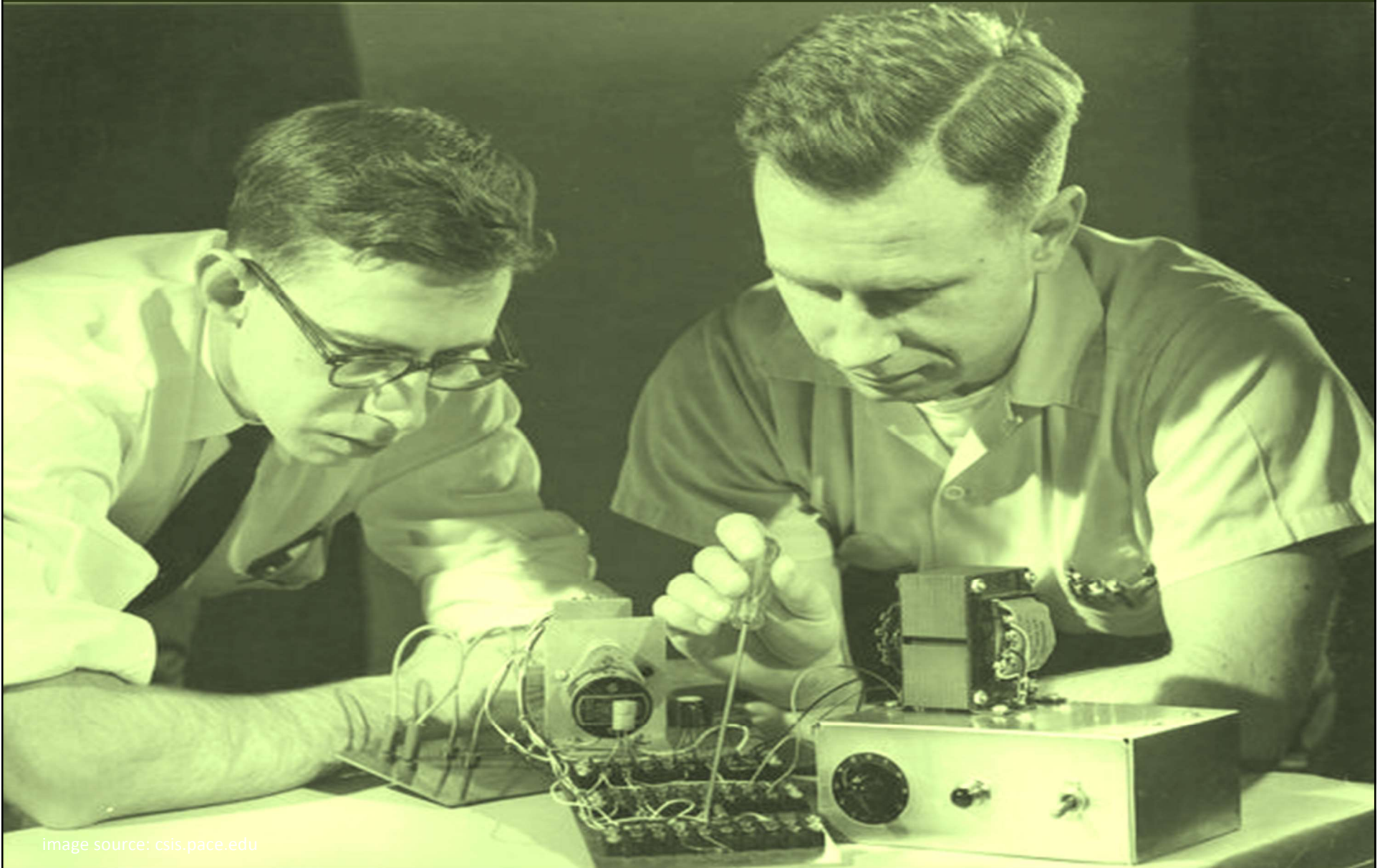
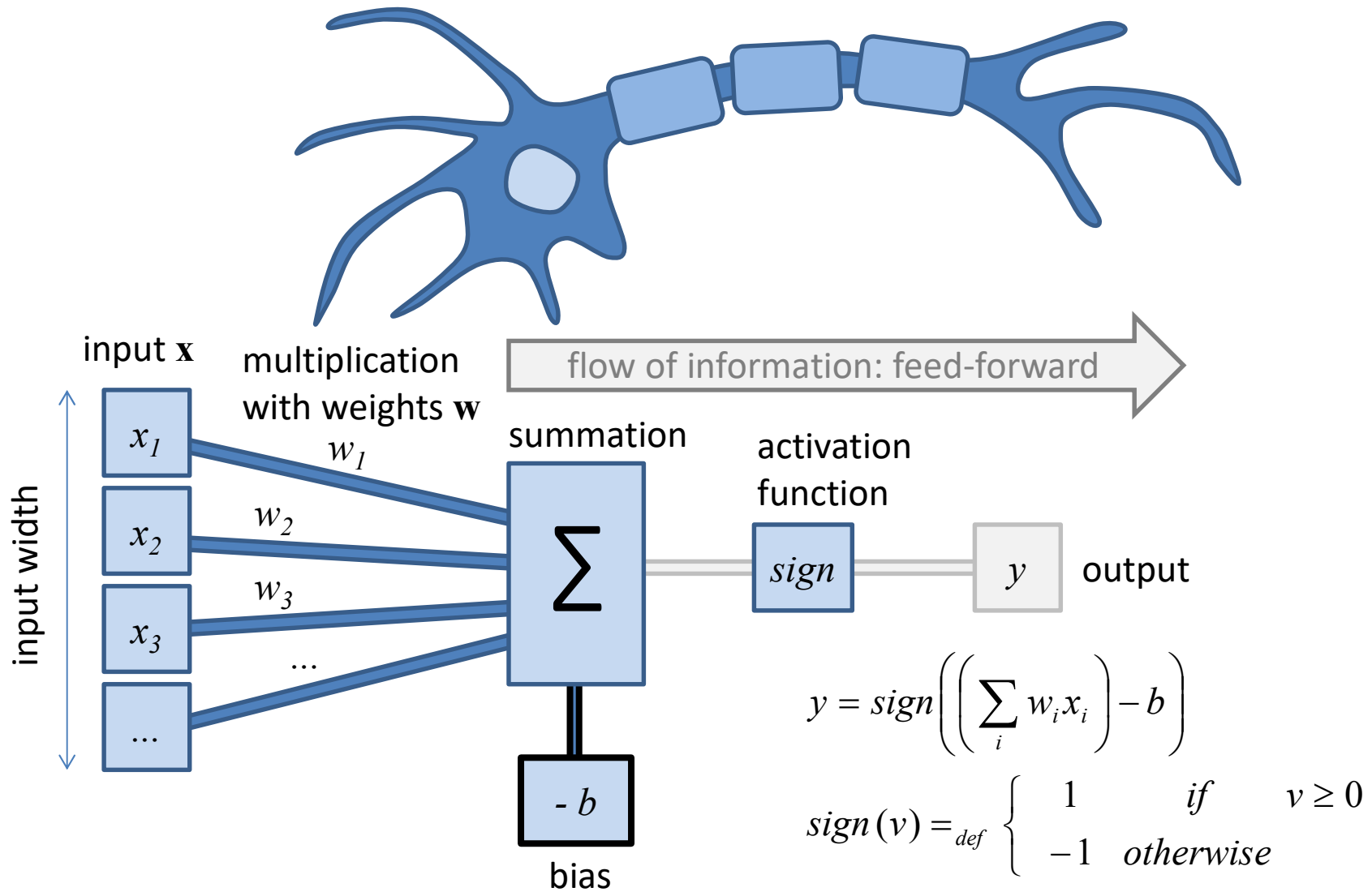
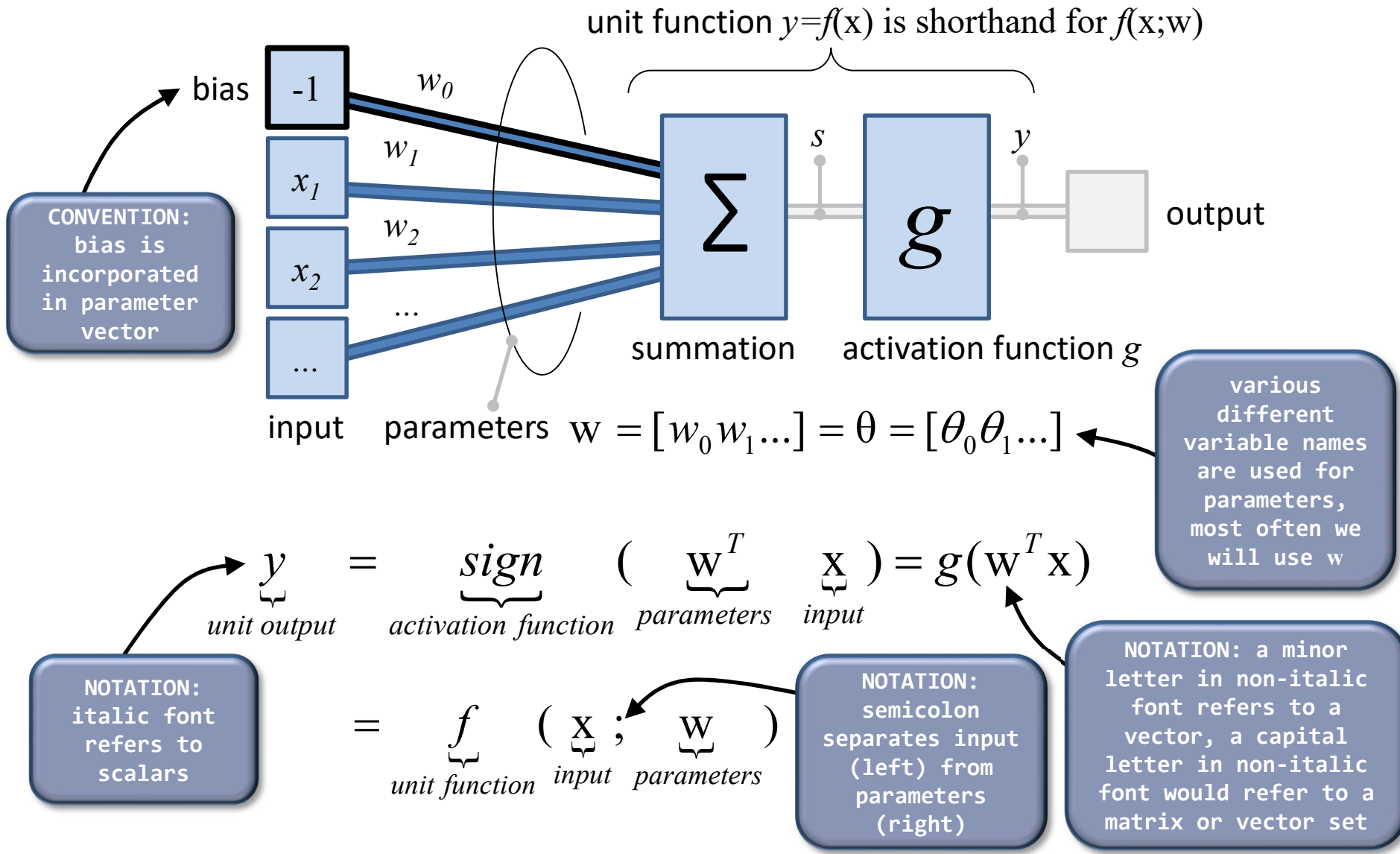


image source: csis.pace.edu

Simplification of a Neuron to a Computational Unit

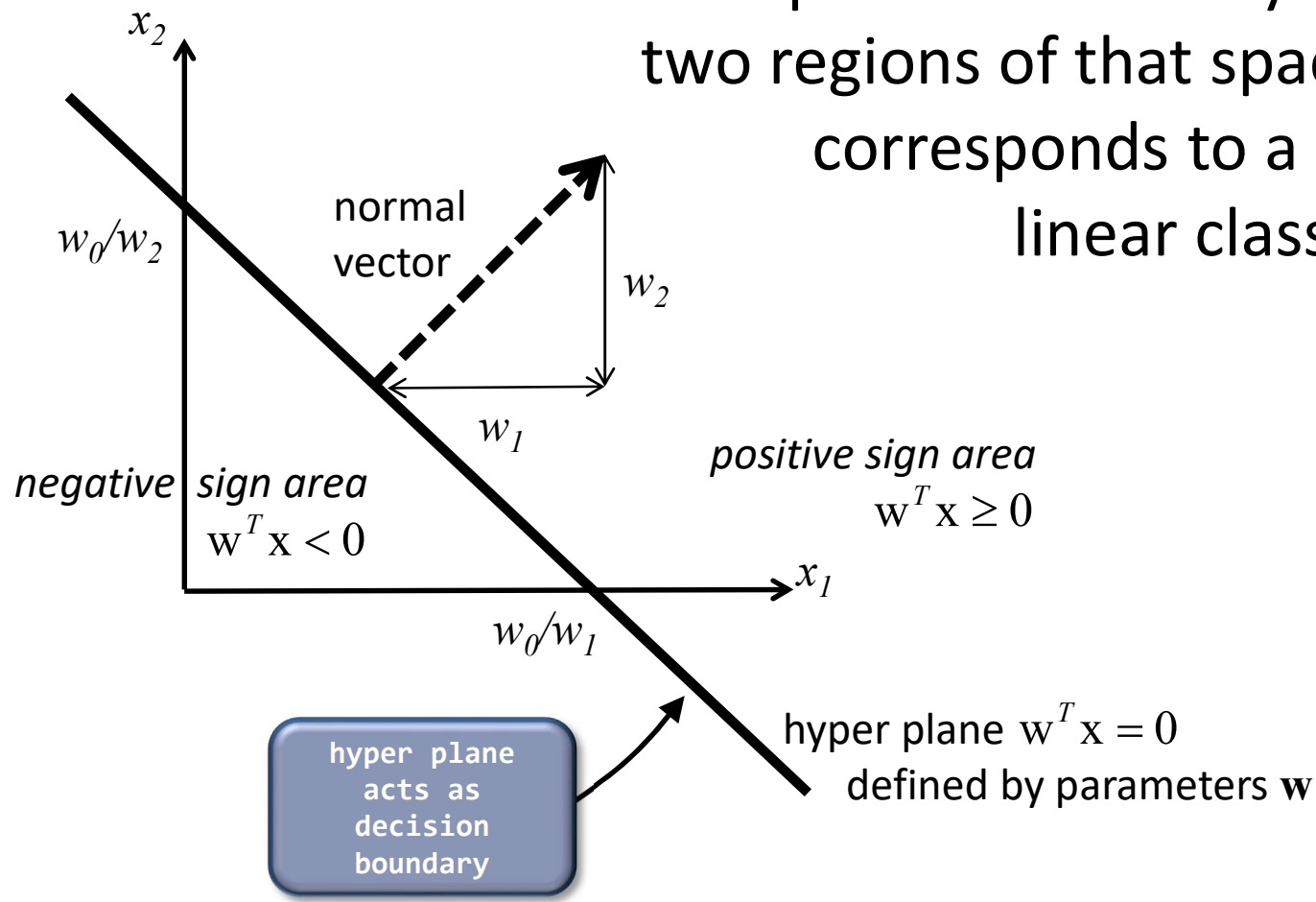


Notational Details for the Perceptron



Geometrical Interpretation of the State Space

The basic Perceptron defines a hyper plane $0 = \mathbf{w}^T \mathbf{x}$ in x -state space that linearly separates two regions of that space (which corresponds to a two-class linear classification)



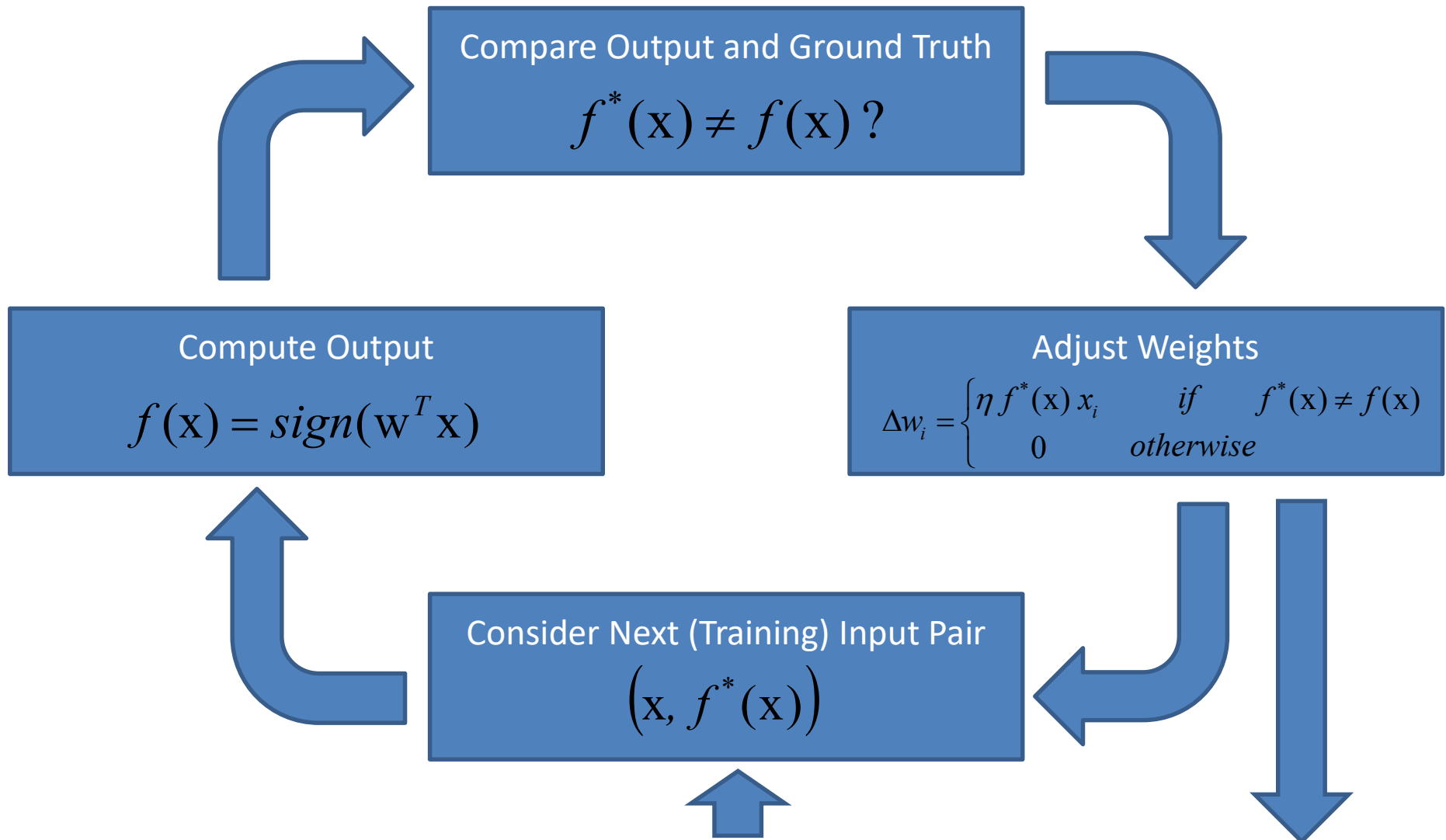
Basic Perceptron (Supervised) Learning Rule

- **Idea:** whenever the system produces a misclassification with current weights, adjust weights by $\Delta \mathbf{w}$ towards a better performing weight vector:

$$\underbrace{\Delta \mathbf{w}}_{\text{update}} = \begin{cases} \eta f^*(\mathbf{x}) \mathbf{x} & \text{if } \overbrace{f^*(\mathbf{x})}^{\text{ground truth}} \neq \overbrace{f(\mathbf{x})}^{\text{actual output}} \\ 0 & \text{otherwise} \end{cases}$$

... where η is the learning rate.

Training a Single-Layer Perceptron



Perceptron Learning Example: OR

Perceptron Training Attempt of **OR** using

$$\Delta w = \eta (f^*(\mathbf{x}) - f(\mathbf{x})) \mathbf{x}; \quad \eta = 0.5$$

OR

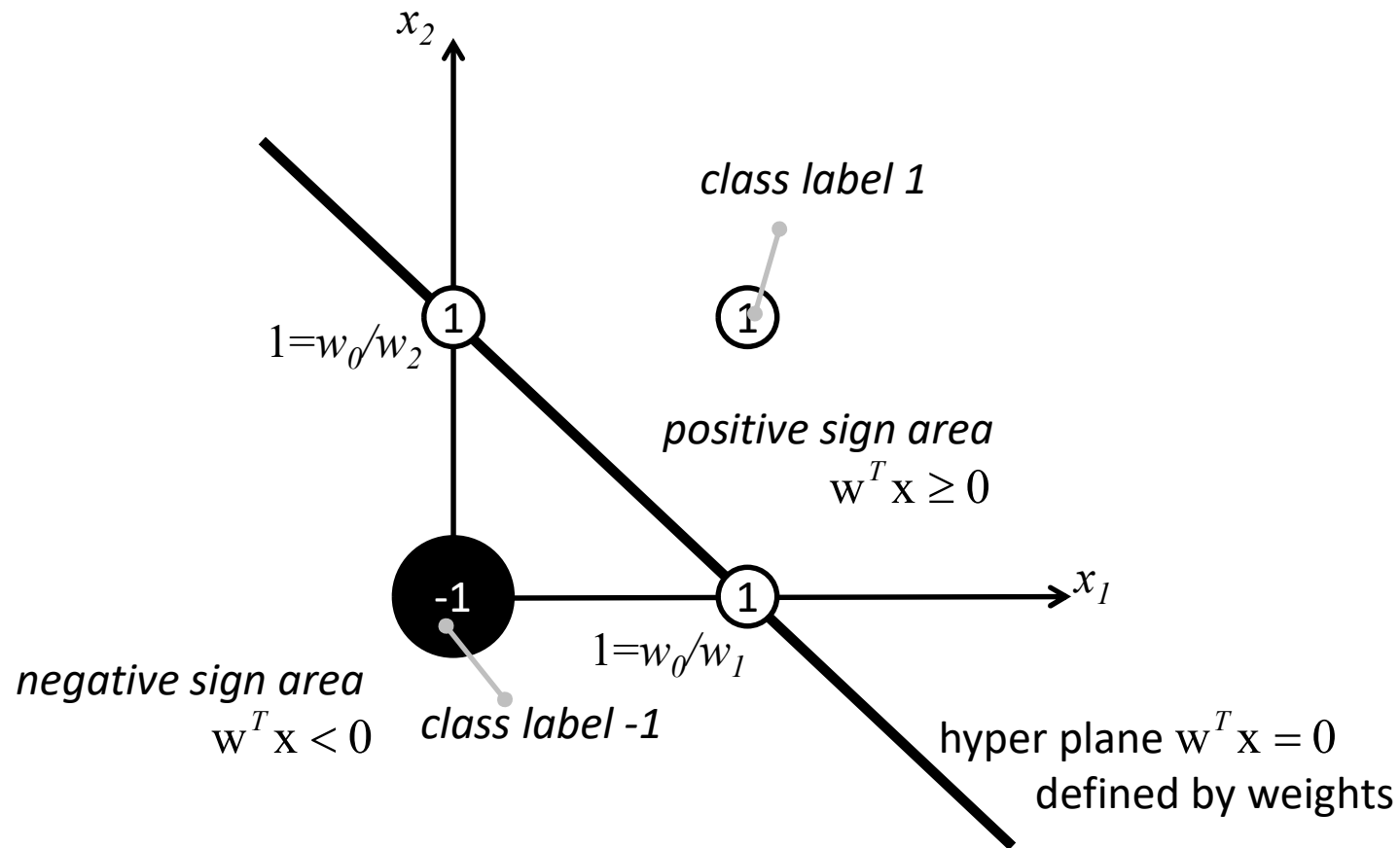
learning progress sampling some (\mathbf{x}, f^*)

x_0	x_1	x_2	parameters w	f	f^*	update Δw
-1	0	0	(0,0,0)	1	-1	(1,0,0)
-1	1	0	(1,0,0)	-1	1	(-1,1,0)
-1	0	0	(0,1,0)	1	-1	(1,0,0)
-1	0	1	(1,1,0)	-1	1	(-1,0,1)
-1	0	0	(0,1,1)	1	-1	(1,0,0)
-1	0	1	(1,1,1)	1	1	(0,0,0)
-1	1	0	(1,1,1)	1	1	(0,0,0)
-1	1	1	(1,1,1)	1	1	(0,0,0)
-1	0	0	(1,1,1)	-1	-1	(0,0,0)
...

x_1	x_2	f^*
0	0	-1
0	1	1
1	0	1
1	1	1

encoding could be changed to traditional value 0 by adjusting the output of the sign function to 0; training algorithm still valid

Geometrical Interpretation of OR Space Learned



Larger Example Visualisation

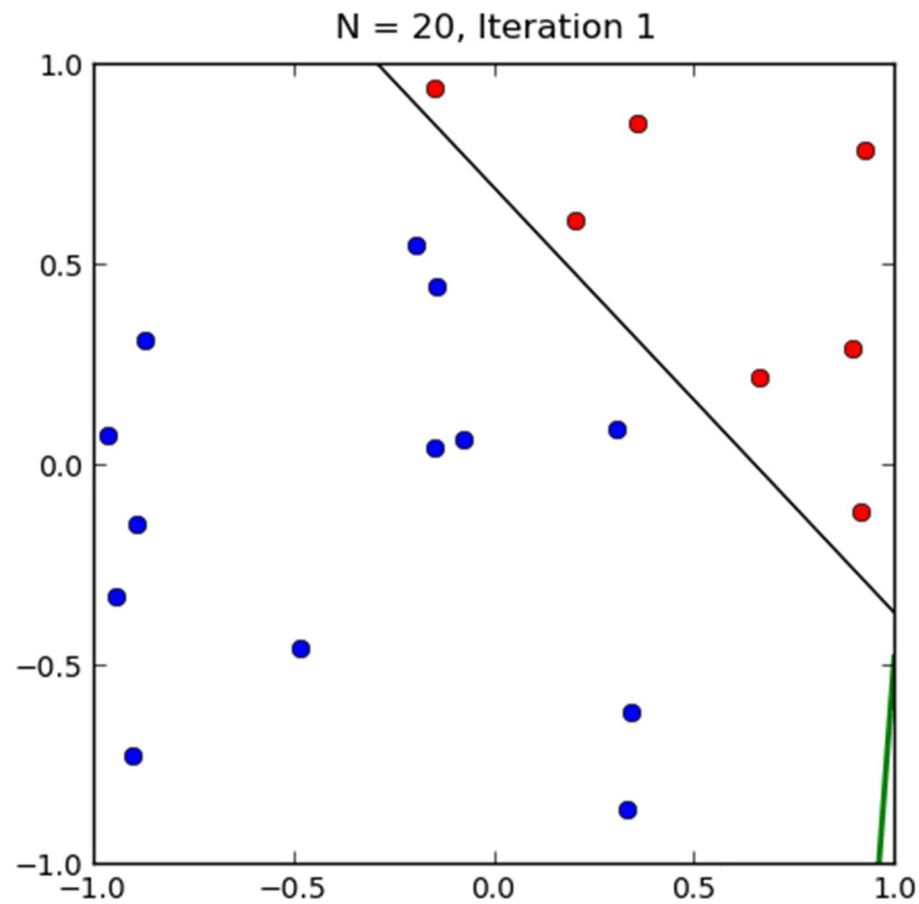


image source: datasciencelab.wordpress.com

COST FUNCTIONS



Cost (or Loss) Functions

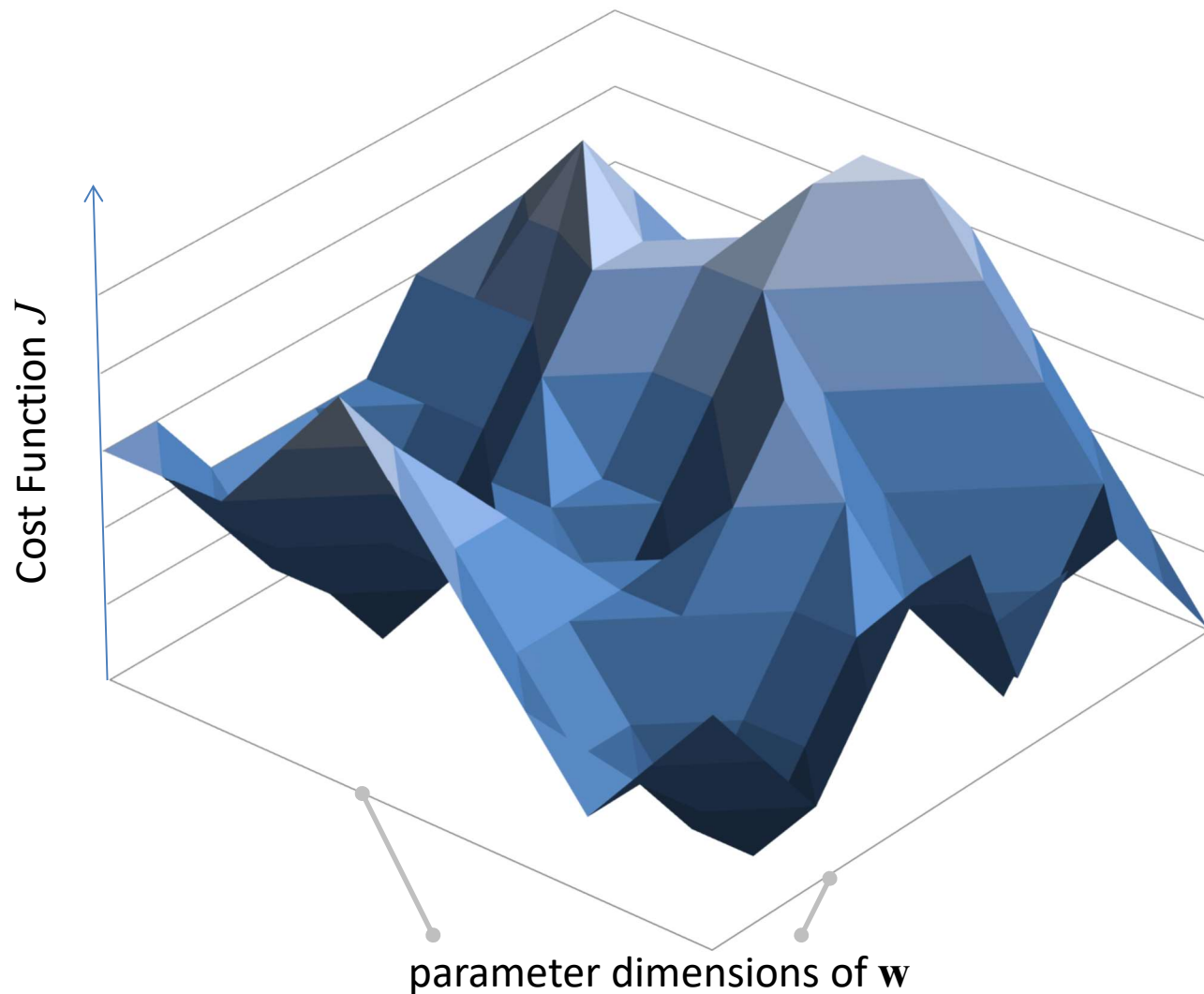
Idea: Given a set \mathbf{X} of input vectors \mathbf{x} of one or more variables and a parameterisation \mathbf{w} , a Cost Function is a map J onto a real number representing a cost or loss associated with the input configurations.
(Negatively related to ‘goodness of fit’.)

Expected Loss: $J(\mathbf{X}; \mathbf{w}) = \mathbb{E}_{(\mathbf{x}, f^*(\mathbf{x})) \sim p} L(f(\mathbf{x}; \mathbf{w}), f^*(\mathbf{x}))$

Empirical Risk: $J(\mathbf{X}; \mathbf{w}) = \frac{1}{|\mathbf{X}|} \sum_{\mathbf{x} \in \mathbf{X}} L(f(\mathbf{x}; \mathbf{w}), f^*(\mathbf{x}))$

MSE Example: $MSE_{loss} = \underbrace{J(\mathbf{X}; \mathbf{w})}_{\text{loss function}} = \frac{1}{|\mathbf{X}|} \sum_{\mathbf{x} \in \mathbf{X}} \underbrace{(f(\mathbf{x}; \mathbf{w}) - f^*(\mathbf{x}))^2}_{\text{per-example loss function}}$

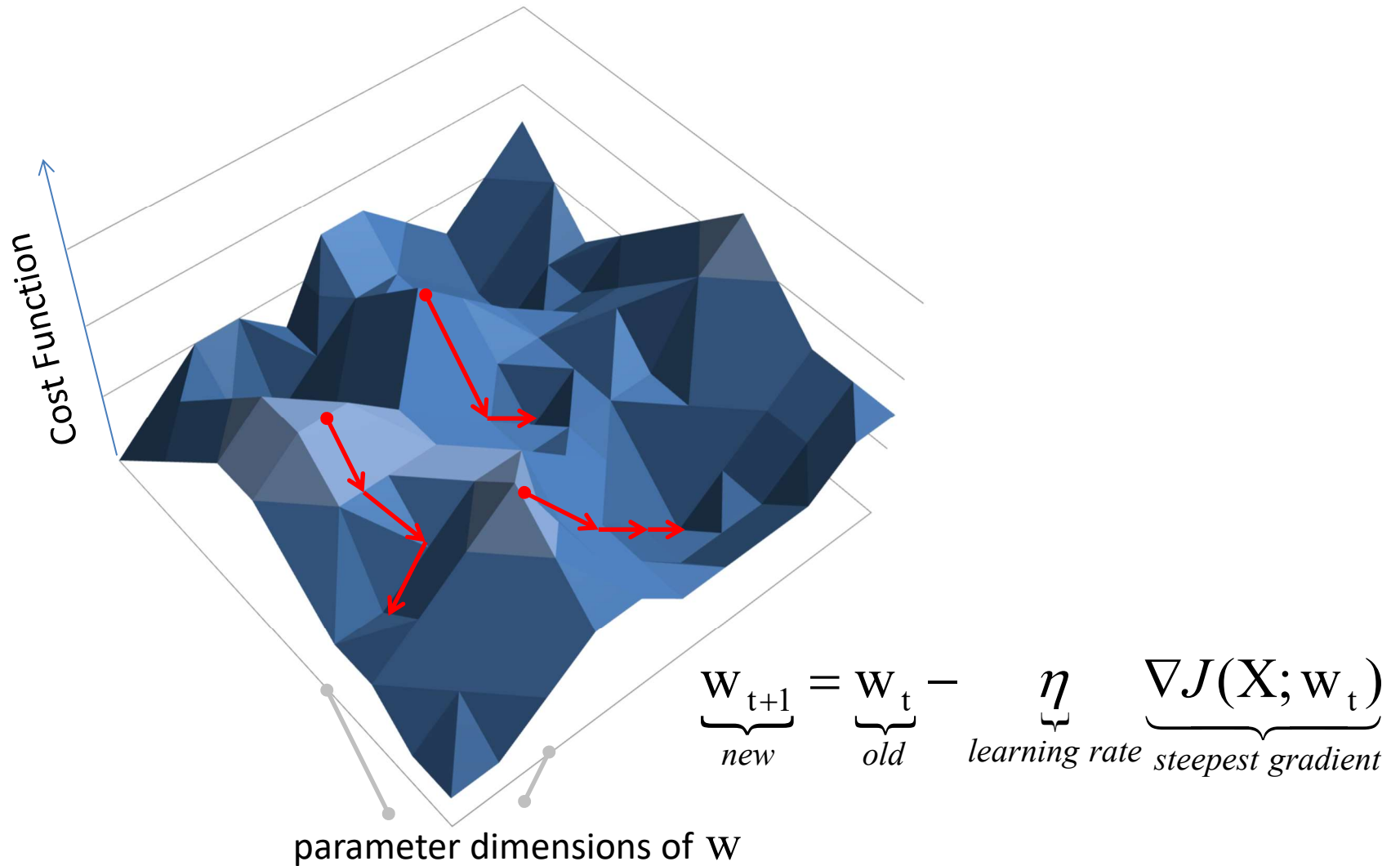
Energy Landscapes over Parameter Space



STEEPEST GRADIENT DESCENT



Idea of 'Steepest' Gradient Descent



The Delta Rule

$$J(\mathbf{X}; \mathbf{w}) = \frac{1}{2 |\mathbf{X}|} \sum_{\mathbf{x} \in \mathbf{X}} \left(\mathbf{w}^T \mathbf{x} - f^*(\mathbf{x}) \right)^2$$

MSE-type cost function with identity function as activation function

$$\Delta \mathbf{w} = -\eta \nabla J(\mathbf{X}; \mathbf{w})$$

weight vector change is modelled as a move along the steepest descent

change for a single weight w_k

$$\Delta w_k = -\eta \frac{\partial J(\mathbf{X}; \mathbf{w})}{\partial w_k} = -\frac{\eta}{|\mathbf{X}|} \sum_{\mathbf{x} \in \mathbf{X}} x_k \left(\mathbf{w}^T \mathbf{x} - f^*(\mathbf{x}) \right)$$

...and for a single sample...

$$\Delta w_k = -\eta x_k \left(\mathbf{w}^T \mathbf{x} - f^*(\mathbf{x}) \right)$$

this term looks similar to the Perceptron learning rule

$$\Delta \mathbf{w} = -\eta \mathbf{x} \underbrace{\left(\mathbf{w}^T \mathbf{x} - f^*(\mathbf{x}) \right)}_{\delta \text{ is the error derivative}}$$

$$\Delta \mathbf{w} = -\eta \mathbf{x} \delta$$

also known as The Delta Rule (Widrow & Hoff, 1960)

LINEAR SEPARABILITY



Basic Learning Example: XOR

Perceptron Training Attempt of **XOR** using

$$\Delta w = \eta (f^*(x) - f(x)) x; \quad \eta = 0.5$$

XOR

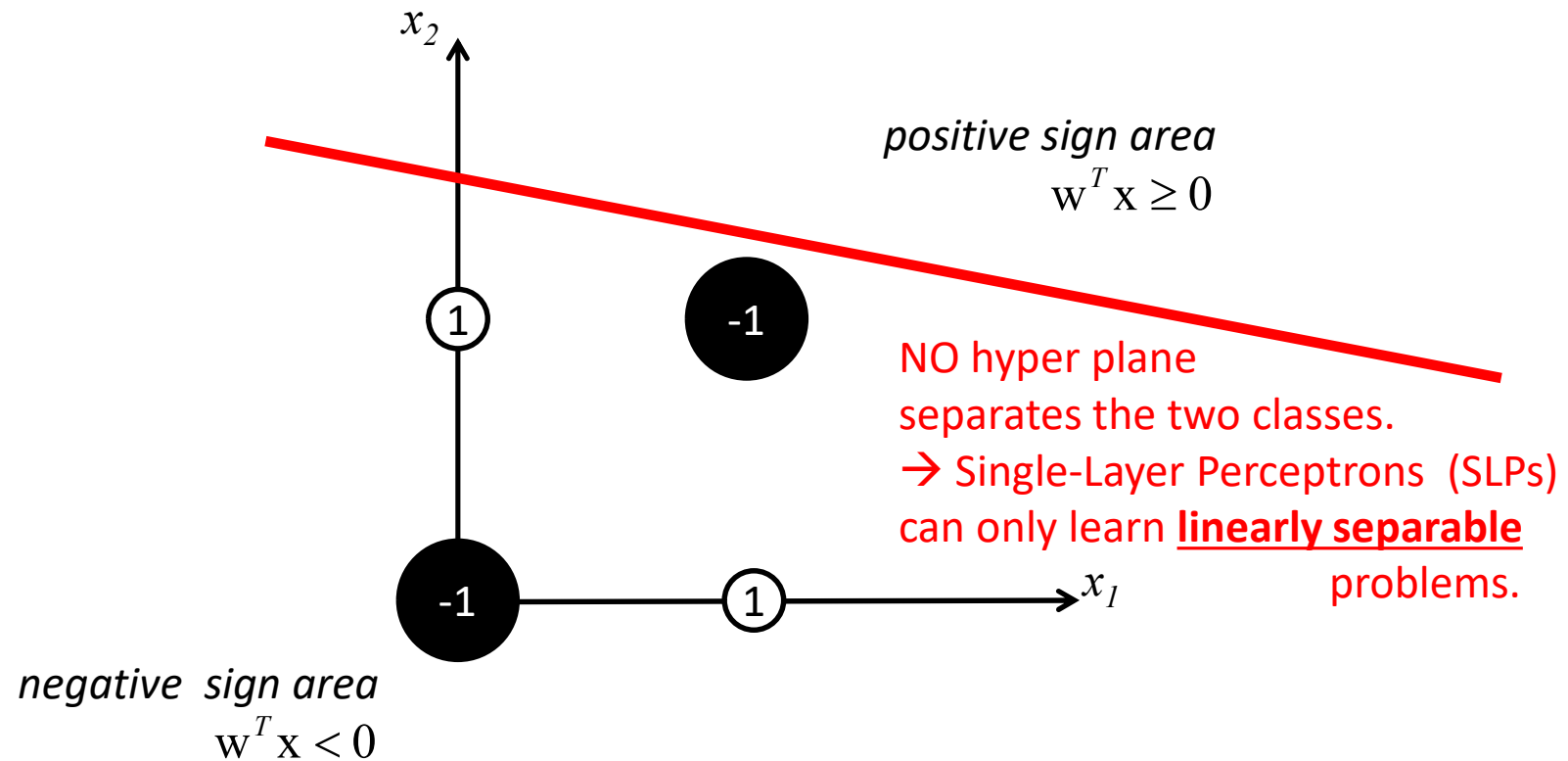
learning progress sampling some (x, f^*)

x_0	x_1	x_2	parameters	f	f^*	update
-1	0	0	(0,0,0)	1	-1	(1,0,0)
-1	1	0	(1,0,0)	-1	1	(-1,1,0)
-1	0	0	(0,1,0)	1	-1	(1,0,0)
-1	0	1	(1,1,0)	-1	1	(-1,0,1)
-1	0	0	(0,1,1)	1	-1	(1,0,0)
-1	0	1	(1,1,1)	1	1	(0,0,0)
-1	1	0	(1,1,1)	1	1	(0,0,0)
-1	1	1	(1,1,1)	1	-1	(1,-1,-1)
-1	1	0	(1,0,0)	-1	1	(-1,1,0)
-1	0	1	(1,1,0)	-1	1	(-1,0,1)
...

x_1	x_2	f^*
0	0	-1
0	1	1
1	0	1
1	1	-1

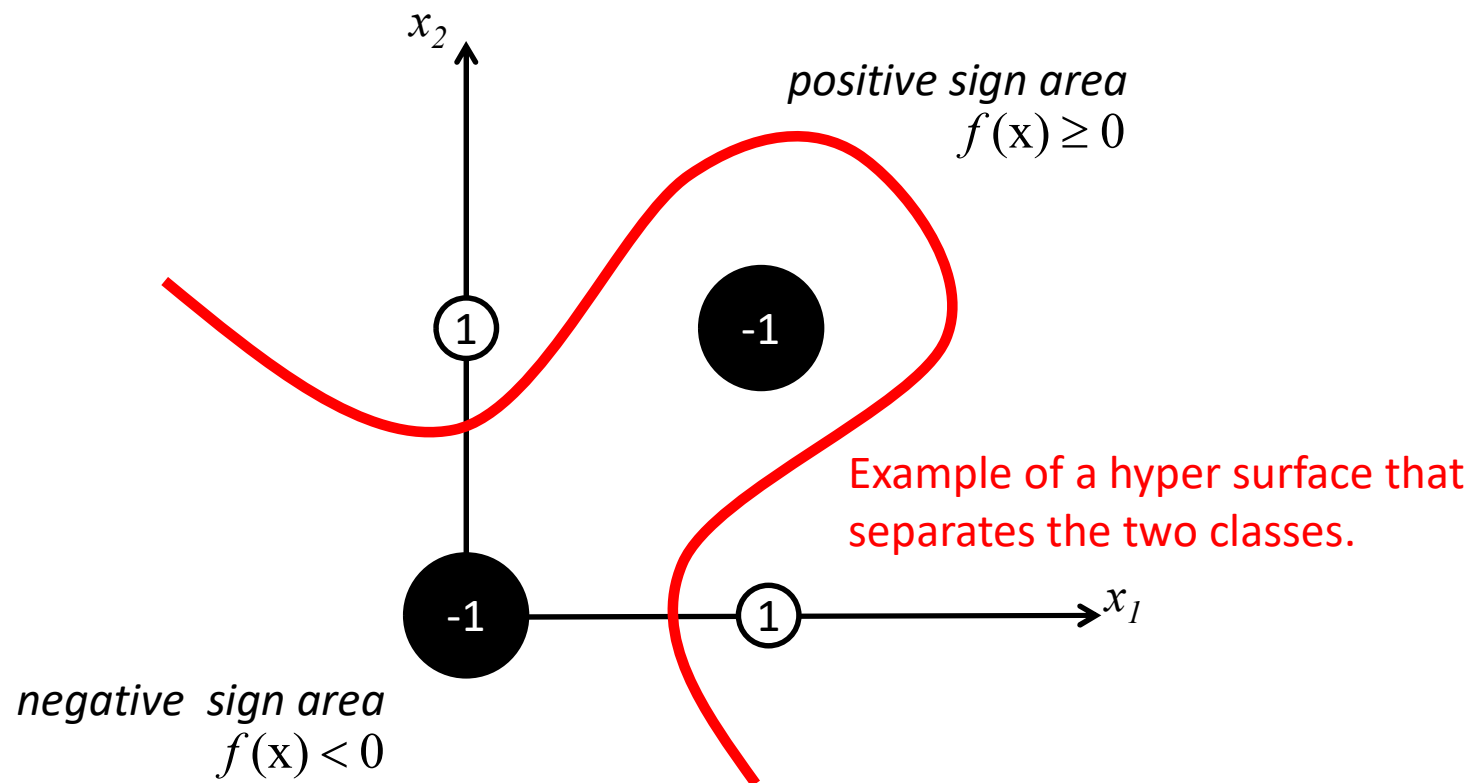
Will the learning process ever produce a solution?

Geometrical Interpretation of XOR Space

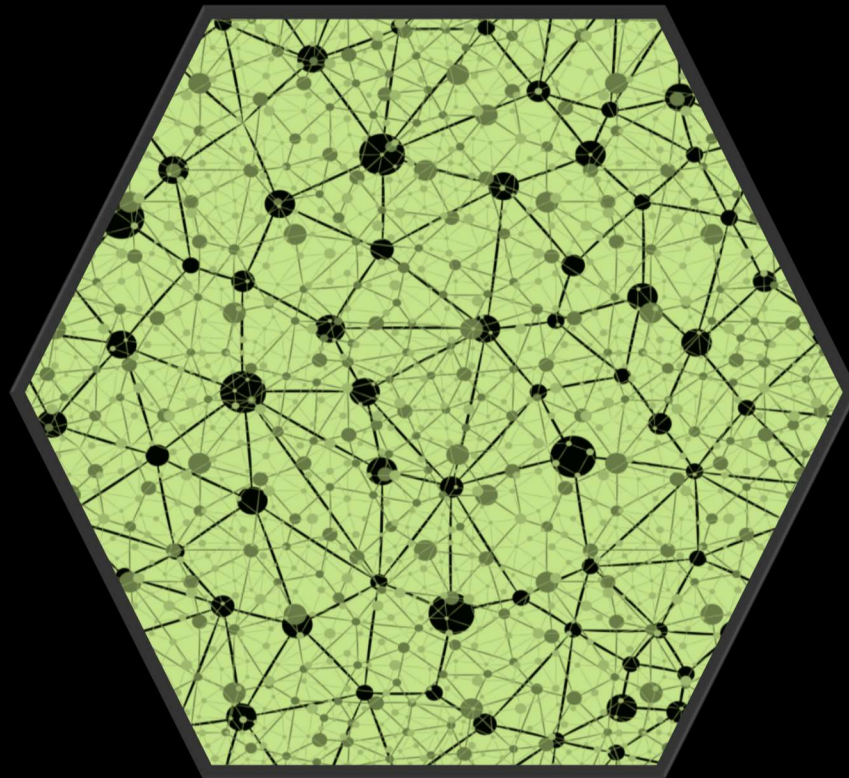


Encoding Arbitrary Decision Boundaries

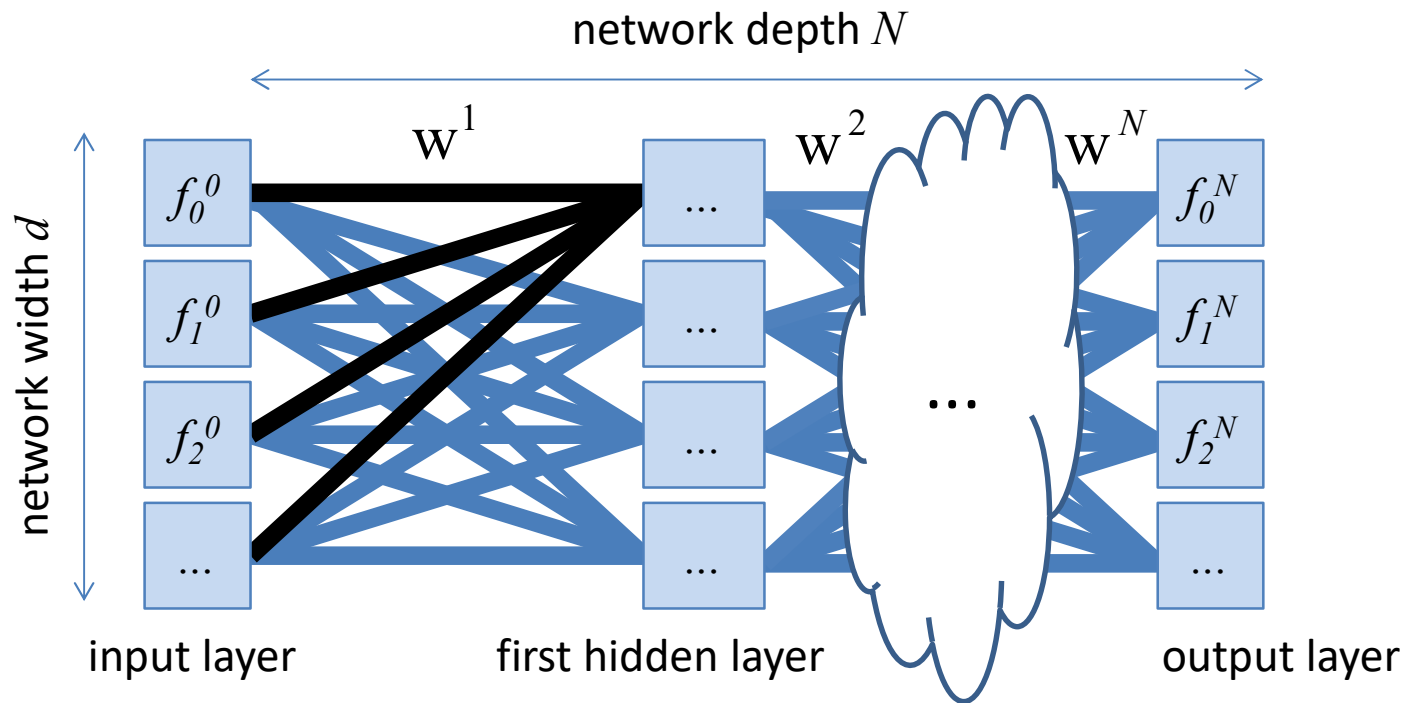
- **Idea:** use of a Multi-Layer Perceptron (MLP) with non-linear activation functions



MULTI-LAYER ARCHITECTURES



Structure and Notation for Deep Architectures



NOTATION: bold math-script used for tensors of order above 2 (basically 3D arrays or higher)

$$\underbrace{\mathbf{f}^N}_{\text{output layer}} =$$

$$\mathbf{f}^{\overset{\text{depth}}{N}} (\dots \mathbf{f}^2(\mathbf{f}^1(\underbrace{\mathbf{f}^0}_{\text{input layer}}; \mathbf{W}^1); \mathbf{W}^2) \dots; \mathbf{W}^N)$$

NOTATION: superscript usually refers to layer number, subscript to position in layer

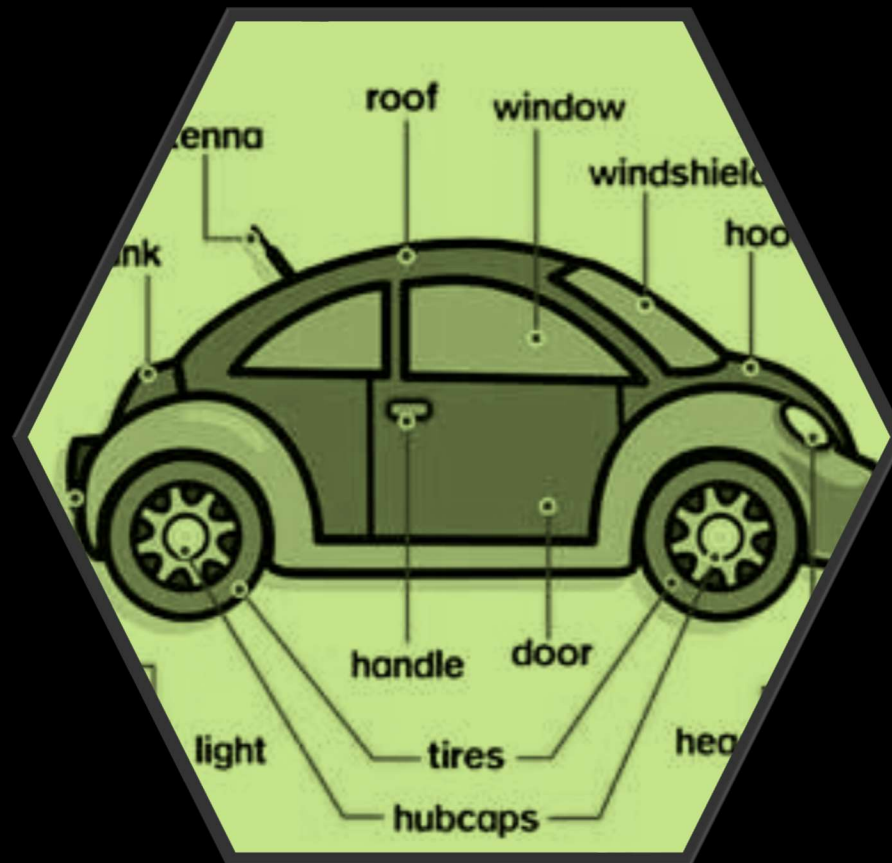
$$\underbrace{\mathbf{W}}_{\text{all parameters}}$$

$$= [\mathbf{W}^1 \mathbf{W}^2 \dots \mathbf{W}^N] = \begin{bmatrix} \dots & \dots & \dots \\ \dots & w_{ij}^l & \dots \\ \dots & \dots & \dots \end{bmatrix}$$

weight w_{ij}^l connects the i^{th} neuron in layer $l-1$ to the j^{th} neuron in layer l

OUTLOOK:

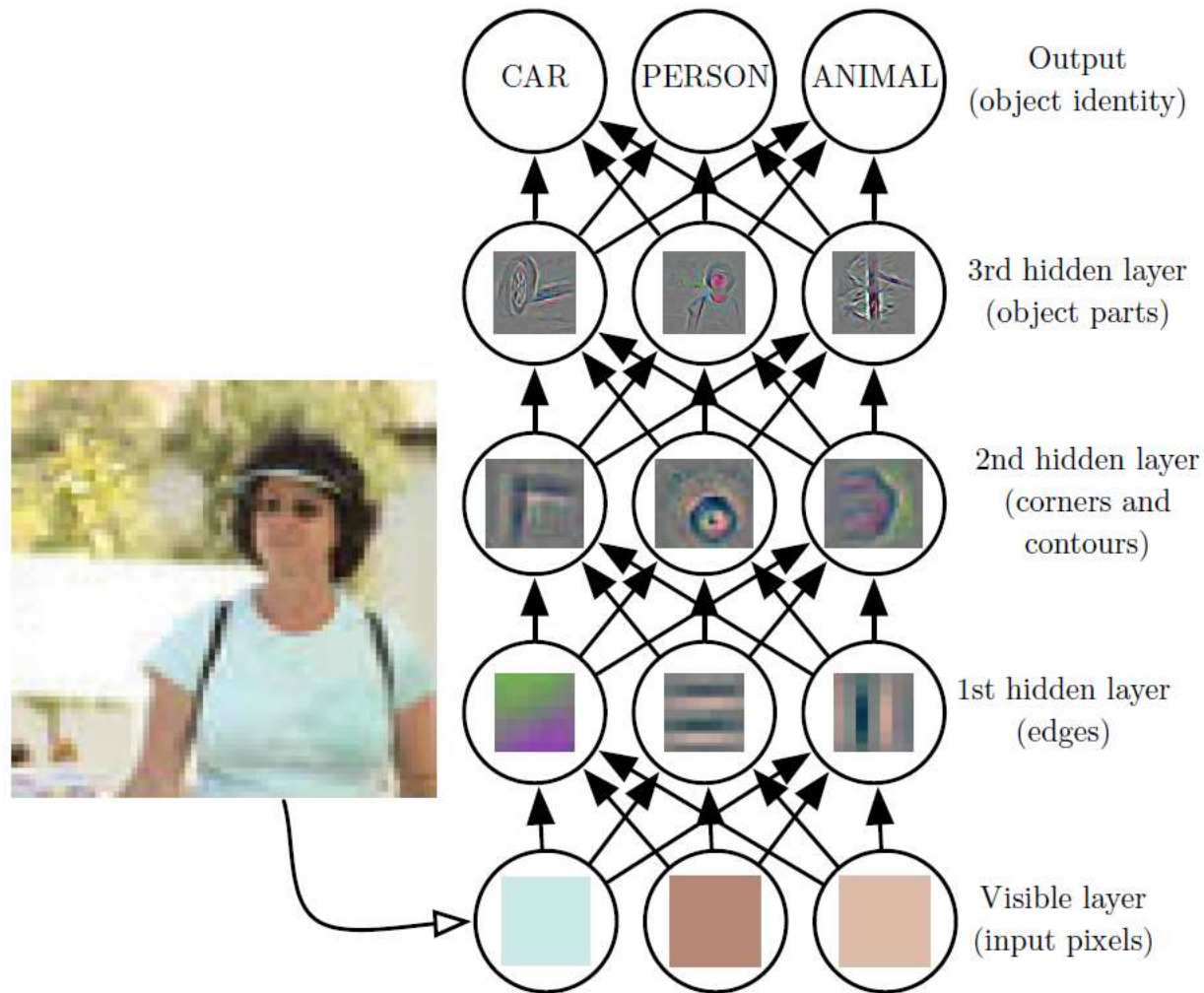
LEARNING REPRESENTATIONS



Representational Power of Feedforward Networks

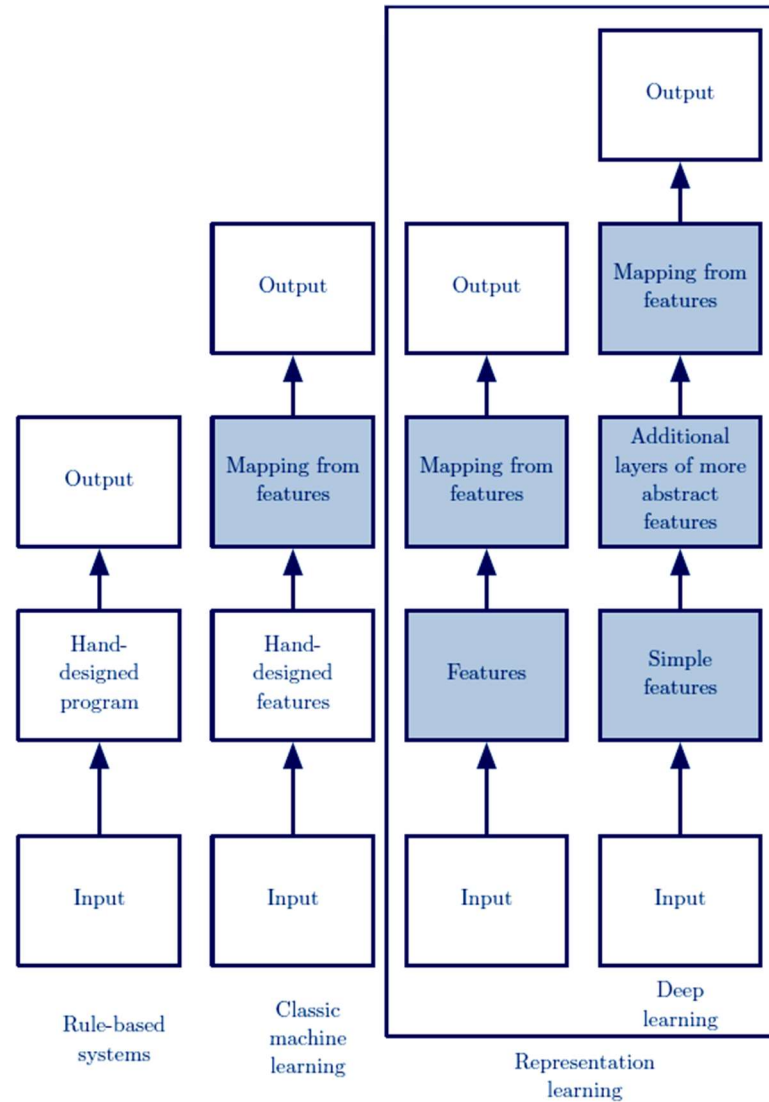
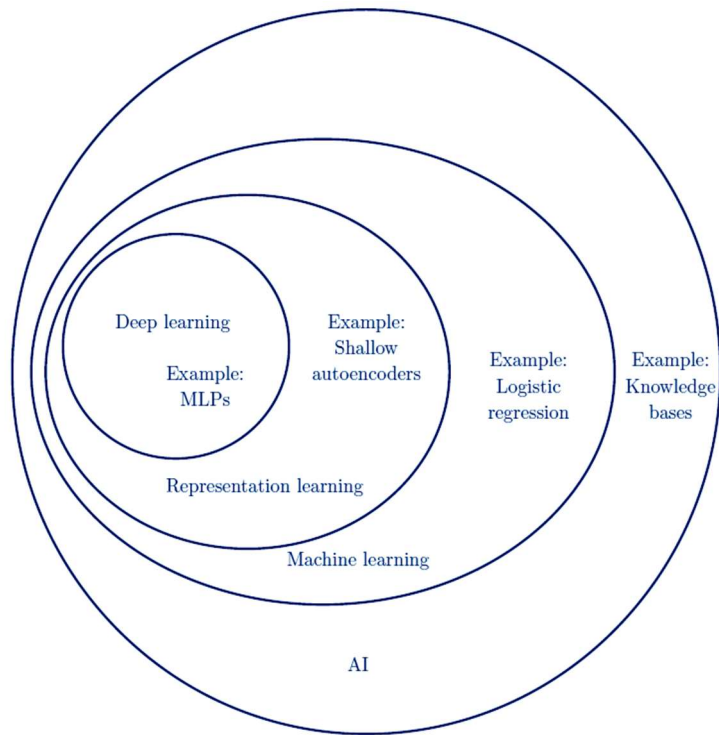
- The basic Perceptron represents a linear classifier.
 - Boolean functions can be represented by layered networks with one hidden layer (networks may be very wide requiring an exponential number of hidden neurons compared to input).
 - Layered networks with one hidden layer can also represent any continuous function [Cybenko 1989; Hornik et al. 1989].
 - Layered networks with two hidden layers can represent any mathematical function [Cybenko 1988].
- long-standing optimism about the potential of neural networks to model learning and intelligent systems
- question arises: why use more than two hidden layers – why is 'deep' advantageous at all? (see Lecture 4)

Deep Composition

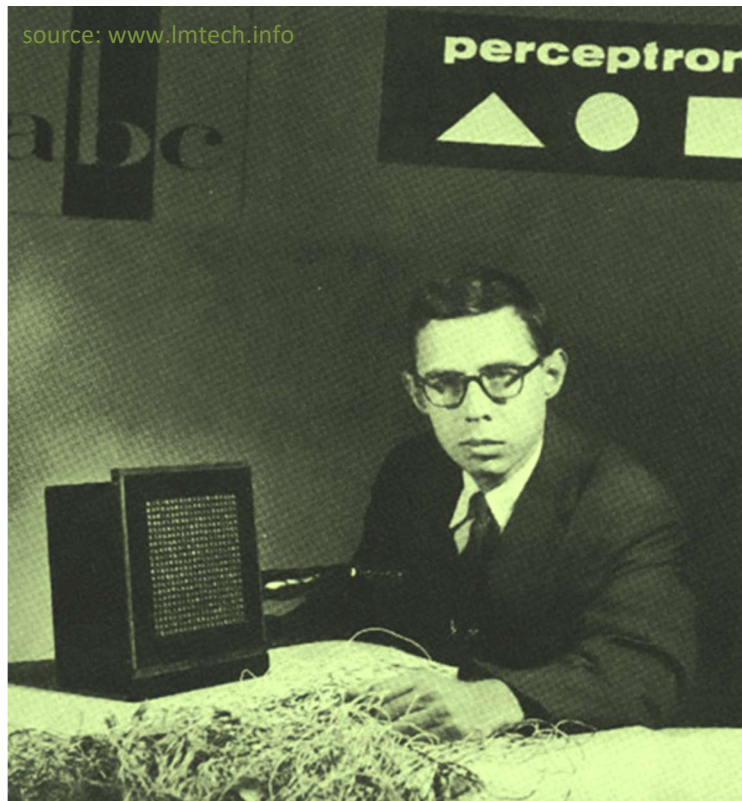


source: Ian Goodfellow, www.deeplearningbook.org

The Concept of Deep Representation Learning



source: Ian Goodfellow, www.deeplearningbook.org



“It is only after much hesitation that the writer has reconciled himself to the addition of the term "neurodynamics" to the list of such recent linguistic artifacts as "cybernetics", "bionics", "autonomics", "biomimesis", "synnoetics", "intelectronics", and "robotics". It is hoped that by selecting a term which more clearly delimits our realm of interest and indicates its relationship to traditional academic disciplines, the underlying motivation of the perceptron program may be more successfully communicated.”

--- *Frank Rosenblatt*

from “Principles of Neurodynamics: Perceptrons and the Theory of Brain Mechanisms”, Spartan Books, 1962

Next Time: Towards Training Deep Architectures

- Computational Graphs
- Reverse Auto-Differentiation

