# Artificial Neural Networks: Lecture 1
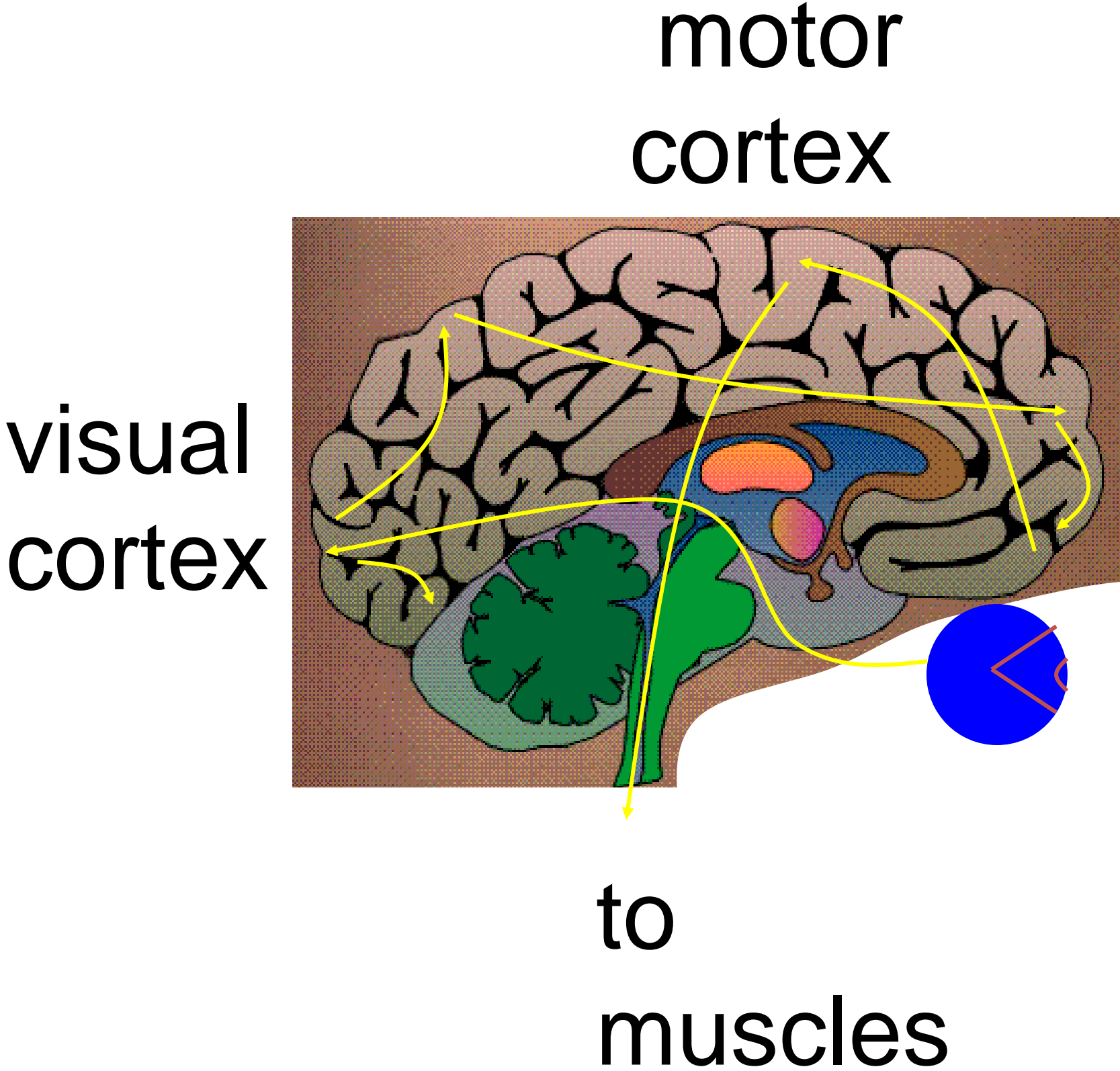## Simple Perceptrons for Classification

Wulfram Gerstner
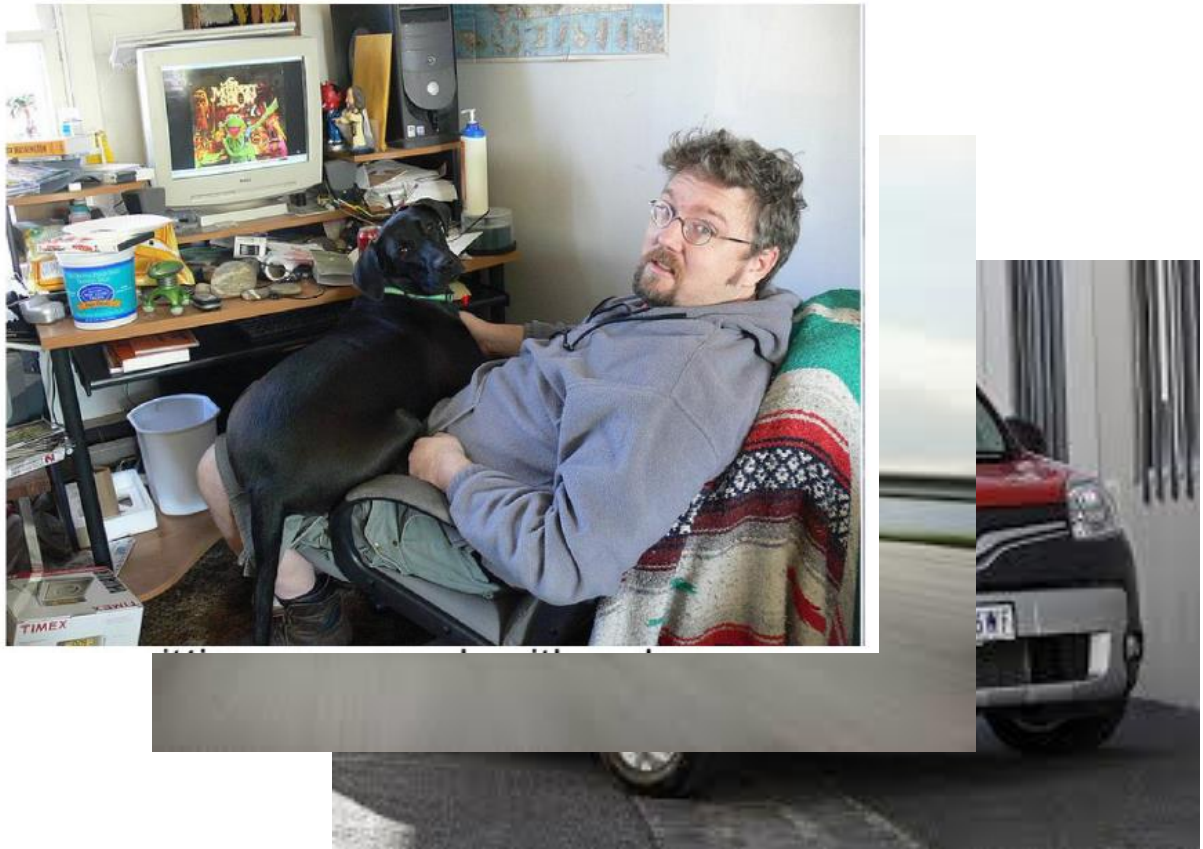EPFL, Lausanne, Switzerland

**Objectives for today:**
- understand classification as a geometrical problem
- discriminant function of classification
- linear versus nonlinear discriminant function
- perceptron algorithm
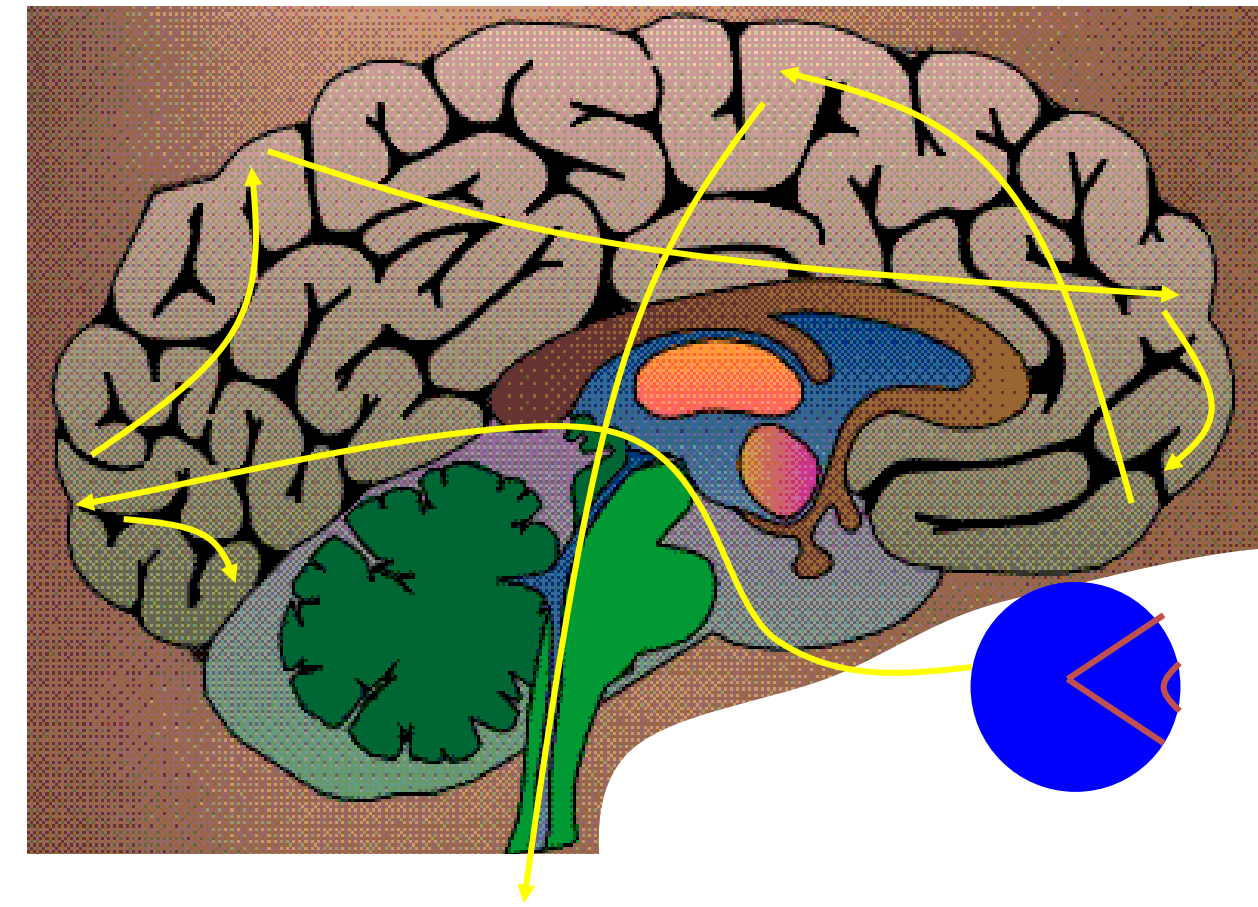- gradient descent for simple perceptrons
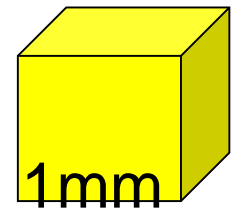
# The brain: Cortical Areas



motor cortex

frontal cortex

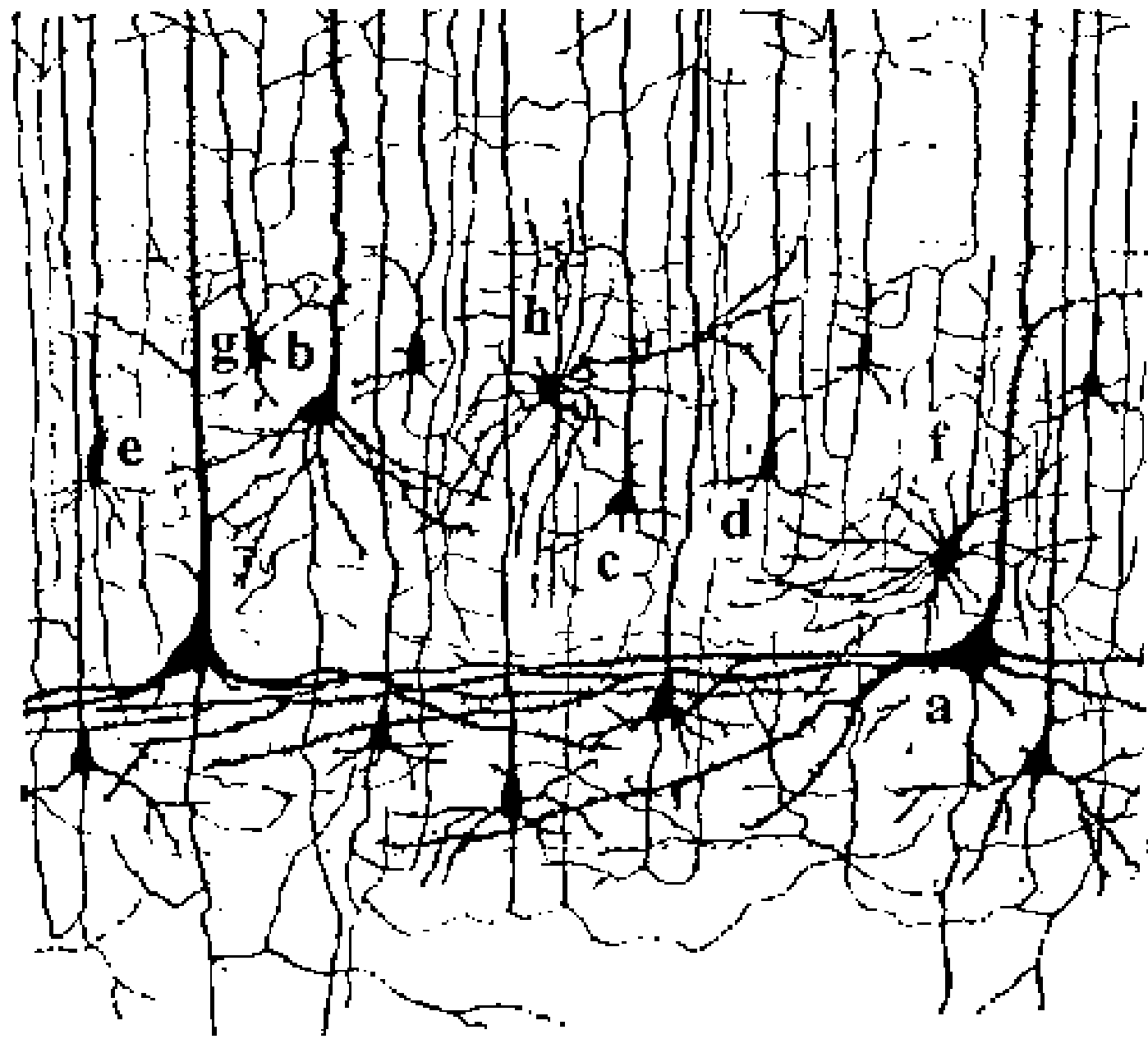visual cortex

to muscles

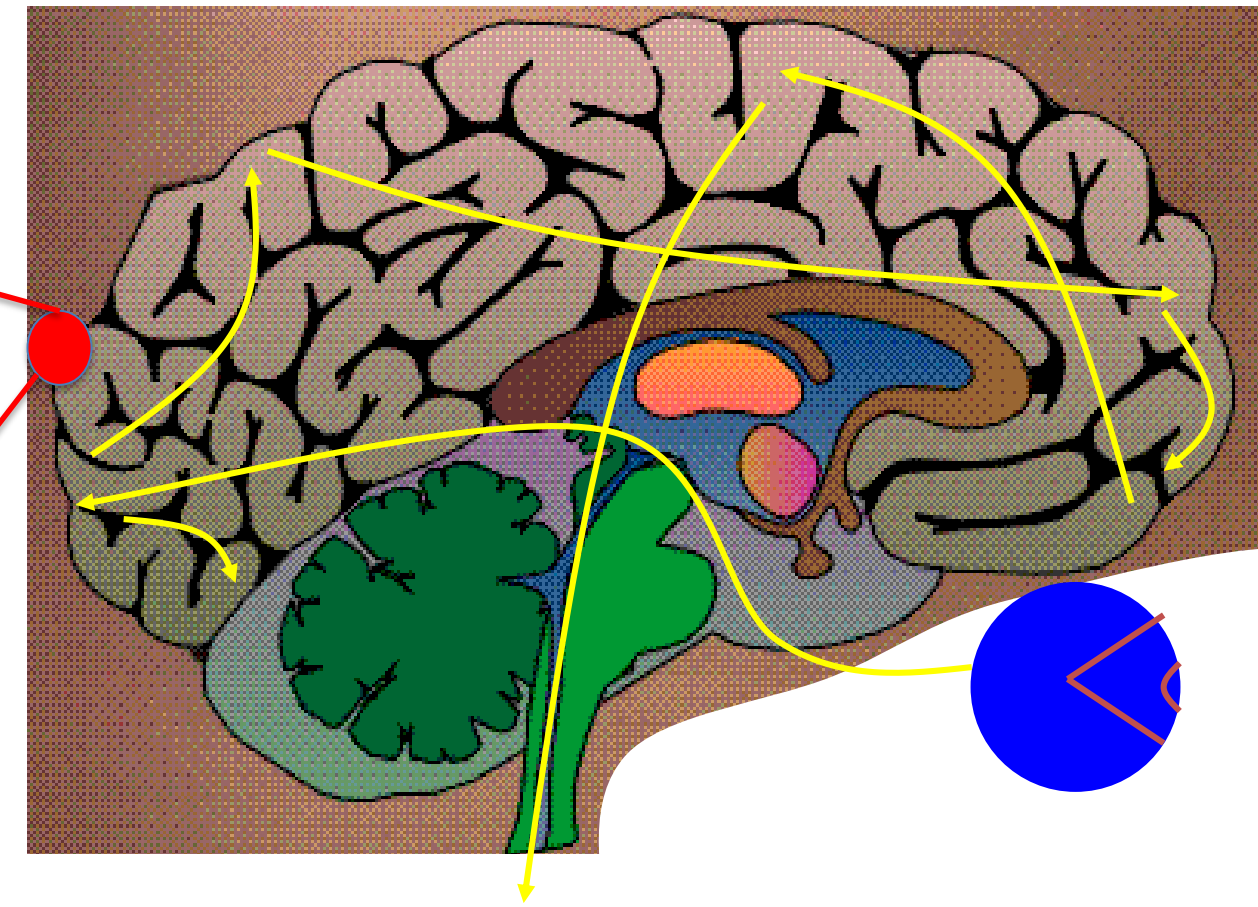# The brain: Cortical Areas

# The Brain: zooming in

1mm

10 000 neurons
3 km of wire
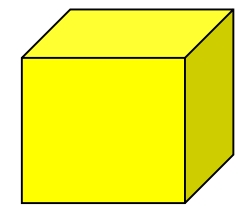
*Ramon y Cajal*

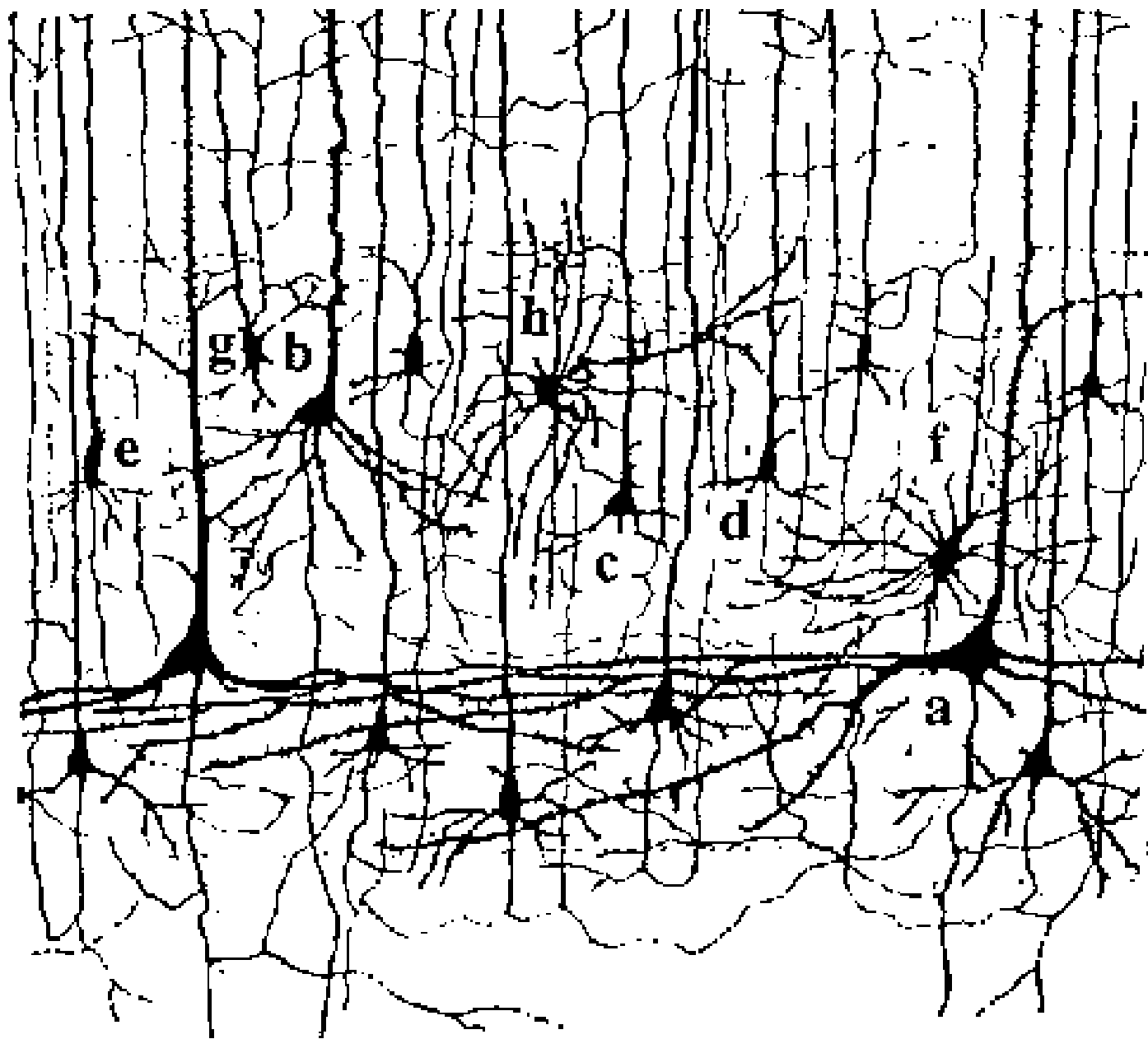# The brain: a network of neurons

1mm


10 000 neurons
3km of wire

Signal:
Action potential (short pulse)


dendrites
soma
axon

*Ramon y Cajal*

# The brain: signal transmission

Signal:
action potential (short pulse)

action potential

More than 1000 inputs

# The brain: neurons sum their inputs



synapse

pulse

$u$

$\vartheta$

$t$

# Summary: the brain is a large network of neurons



🔴 Active neuron

# Learning in the brain: changes between connections

**Neurons**

**Synapse**

**learning = change of connection**

# Artificial Neural Networks

Wulfram Gerstner

EPFL, Lausanne, Switzerland

1. The brain
2. **Artificial Neural Networks**

# Modeling: artificial neurons



pulse

synapse

response

$u$

$\vartheta$

$t$

-responses are added
-pulses created at threshold
-transmitted to other

$\longrightarrow$ Mathematical description

# Modeling: artificial neurons

forget spikes: continuous activity x
forget time: discrete updates

nonlinearity/threshold

activity of output $x_i = g\left(\sum_k w_{ik}\, x_k \atop \vartheta \right)$

$w_{ik}$

weights =
adaptive
parameters

activity of inputs $x_k$

# Neurons and Synapses form a big network



Brain



1mm

10 000 neurons
3 km of wire

10 billions neurons

10 000   connexions/neurons

**memory in the connections**

Distributed Architecture

**No separation of processing and memory**

# Quiz: biological neural networks

[ ] Neurons in the brain have a threshold.
[ ] Learning means a change in the threshold.
[ ] Learning means a change of the connection weights
[ ] The total input to a neuron is the weighted sum of individual inputs
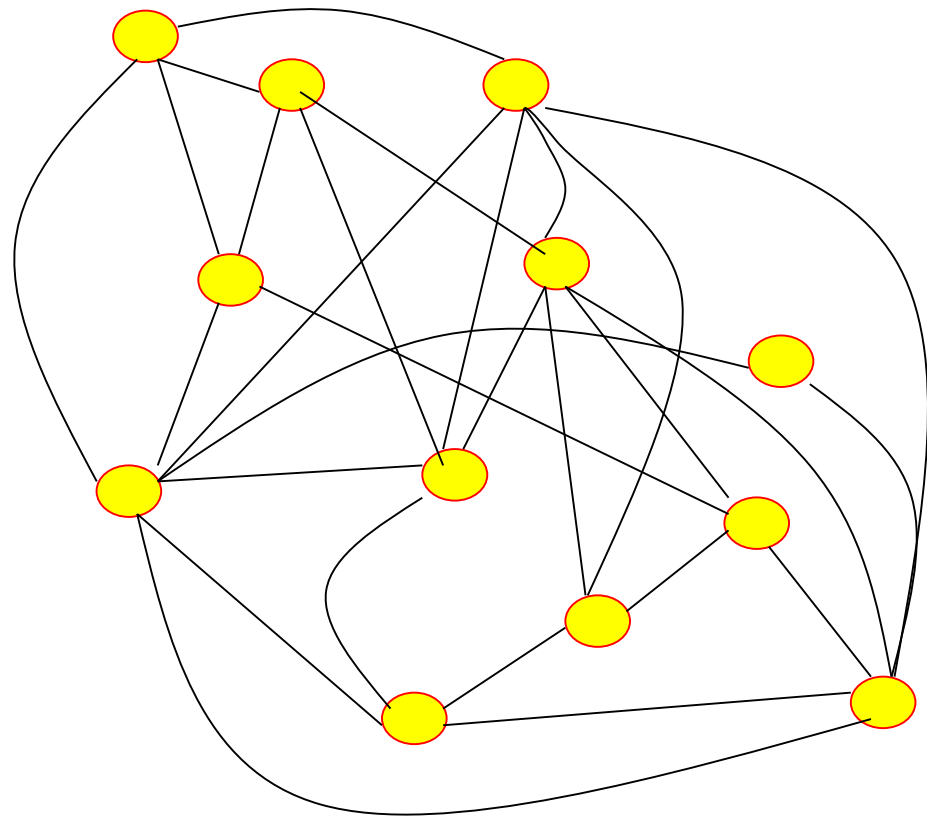[ ] The neuronal network in the brain is feedforward: it has no recurrent connections

# Artificial Neural Networks

Wulfram Gerstner

EPFL, Lausanne, Switzerland

1. The brain
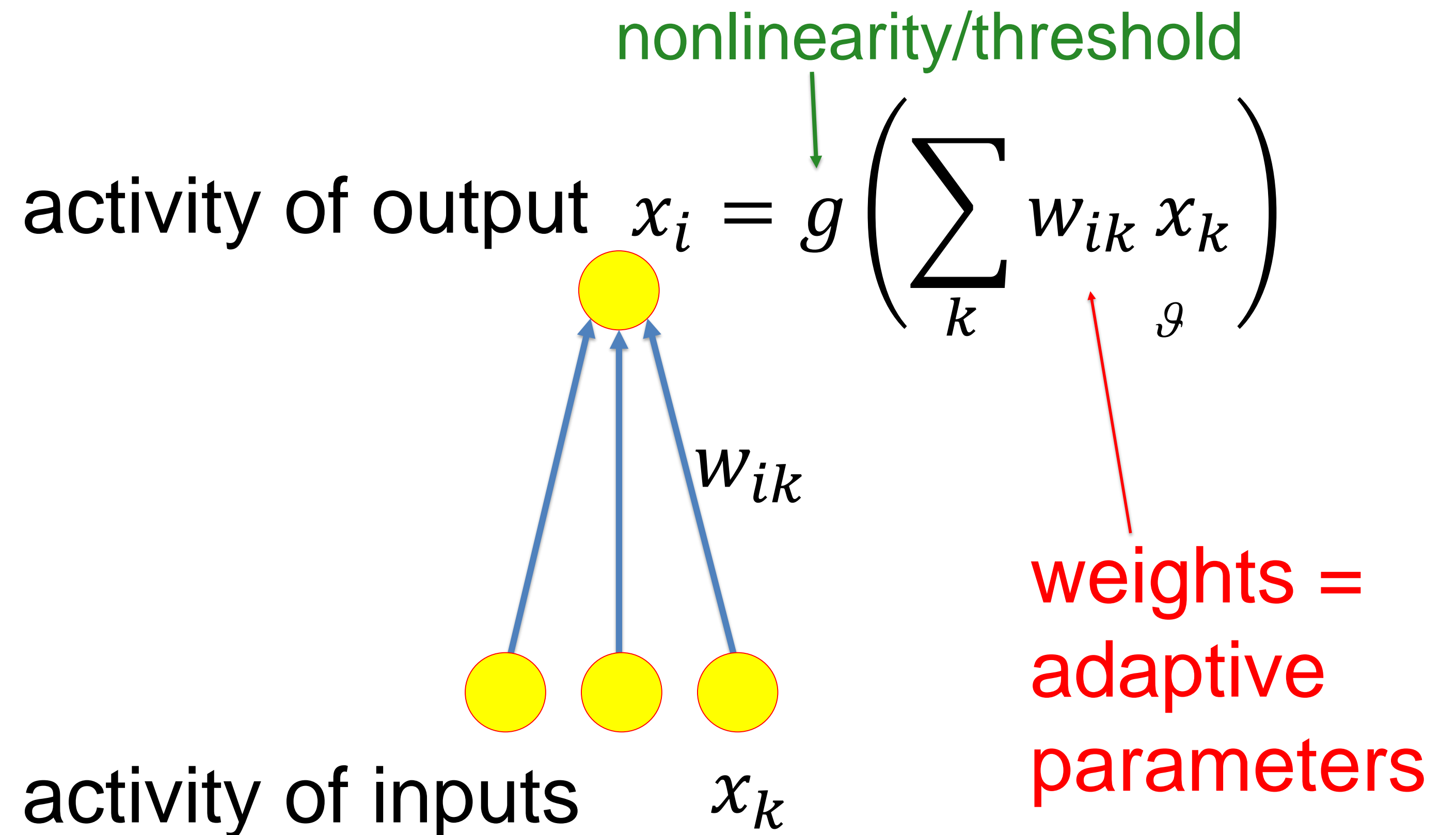2. **Artificial Neural Networks**
   **- artificial neurons**
   **- artificial neural networks for classification**

# Artificial Neural Networks for classification

car

output

feedforward network

input

# Artificial Neural Networks for classification

car    dog

output

input

**Aim of learning:**
Adjust connections such
that output class is correct
(for each input)

# Artificial Neural Networks

Wulfram Gerstner

EPFL, Lausanne, Switzerland

1. The brain
2. Artificial Neural Networks
   - artificial neurons
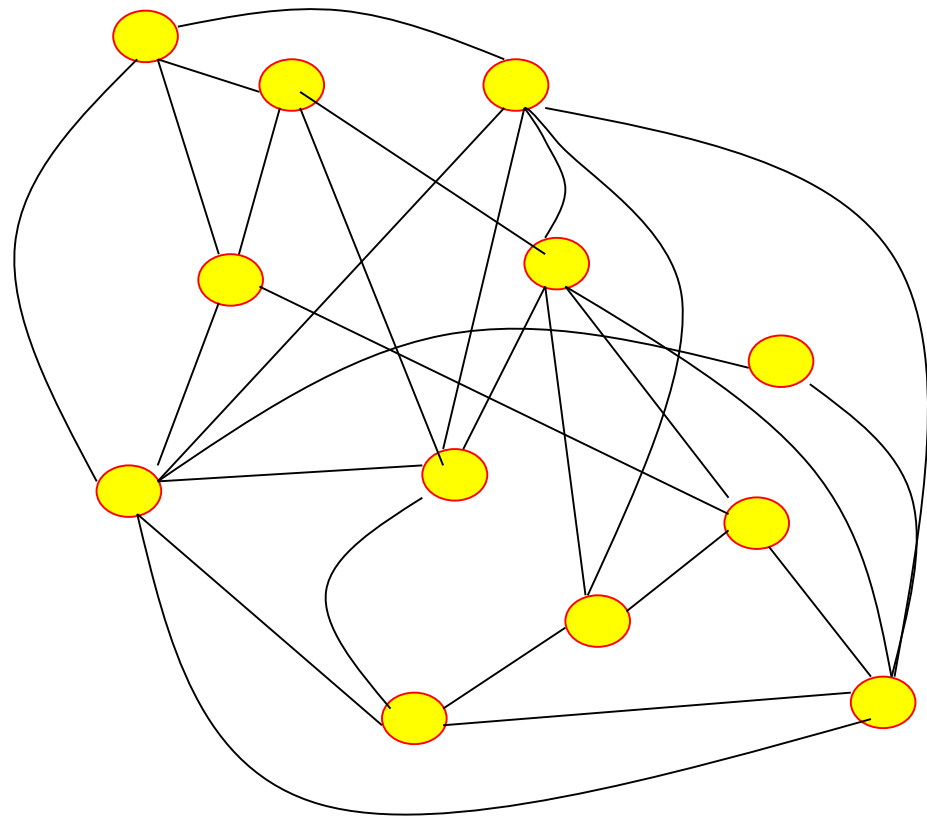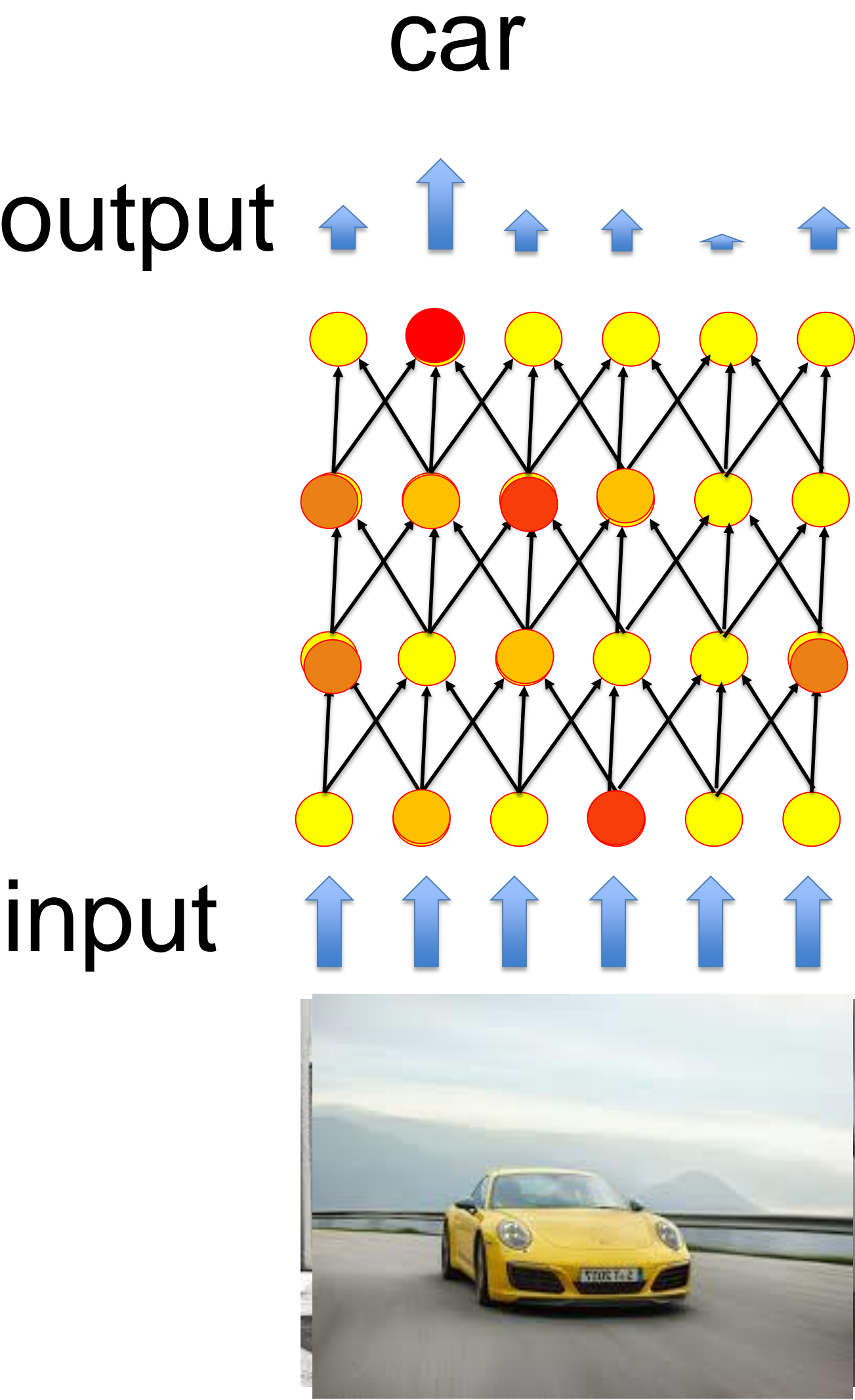   - Neural networks for classification
   - **Neural networks for action learning**

# Deep reinforcement learning

Chess



Go



Artificial neural network (*AlphaZero*) discovers different strategies by playing against itself.

In Go, it beats Lee Sedol

# Reinforcement learning: Learning through rewards (win)

Network for choosing action

action:
*Advance king*

2e output for **value** of action:
*probability to win*

output



input

**learning**:
- change connections

**aim:**
- Predict value of position
- Choose next action to win

# Deep reinforcement learning (alpha zero)

*Silver et al. (2017) , Deep Mind*

output: 4672 actions

**Training** 44Mio games (9 hours)

*advance king*



**Planning: potential sequences** (during 1s before playing next action)

input: 64x6x2x8 neuronss (about 10 000)

# Deep reinforcement learning (alpha zero)

*Silver et al. (2017)*

**Chess:**

-discovers classic openings

-beats best human players

-bets best classic AI algorithms



A10: English Opening
w 20/30/0, b 8/40/2
1...e5 g3 d5 cxd5 ♘f6 ♗g2 ♘xd5 ♘f3

D06: Queens Gambit
w 16/34/0, b 1/47/2
2...c6 ♘c3 ♘f6 ♘f3 a6 g3 c4 a4

A46: Queens Pawn Game
w 24/26/0, b 3/47/0
2...d5 c4 e6 ♘c3 ♗e7 ♗f4 O-O e3
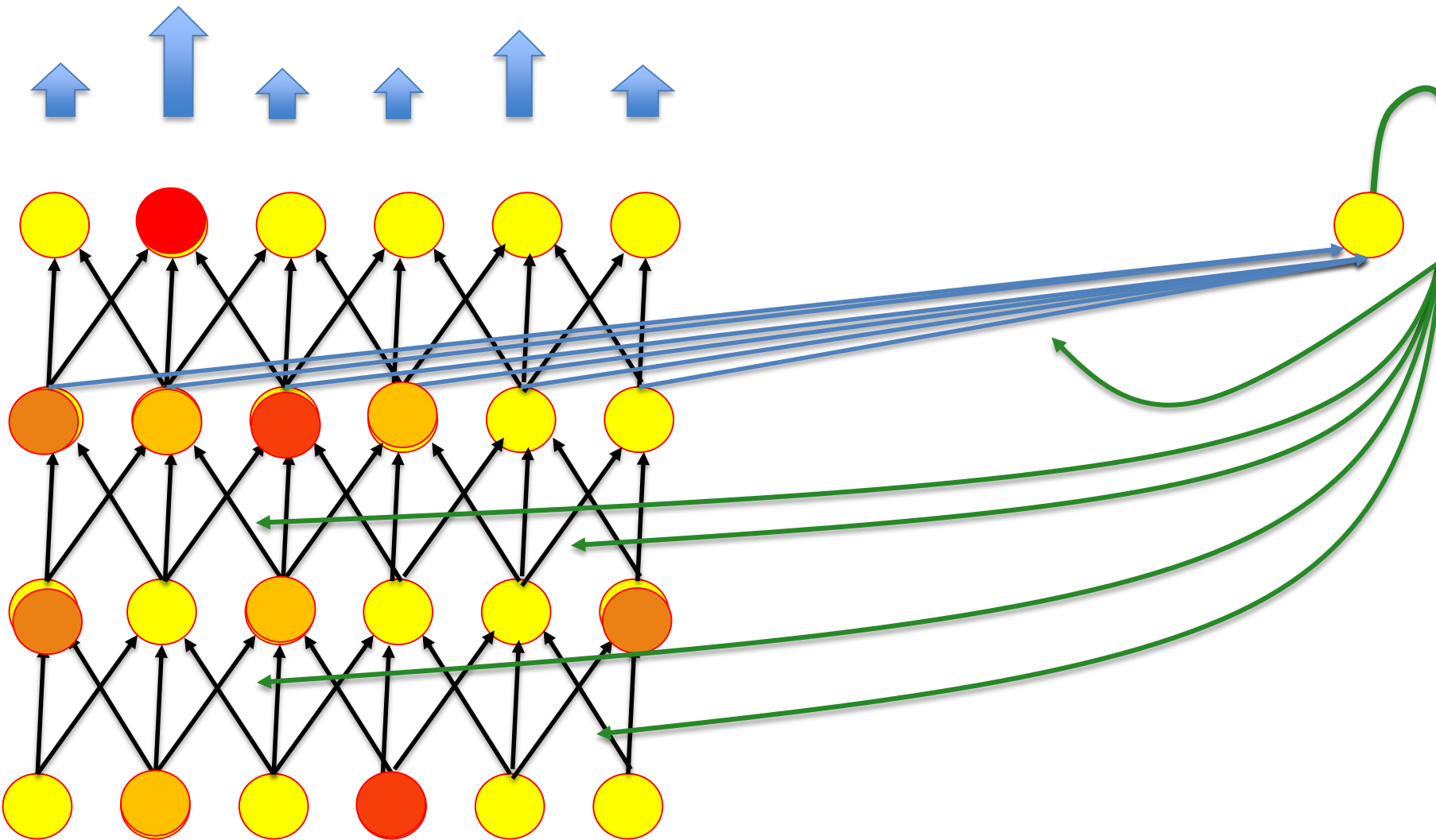
E00: Queens Pawn Game
w 17/33/0, b 5/44/1
3.♘f3 d5 ♘c3 ♗b4 ♗g5 h6 ♕a4 ♘c6

E61: Kings Indian Defence
w 16/34/0, b 0/48/2
3...d5 cxd5 ♘xd5 e4 ♘xc3 bxc3 ♗g7 ♗e3

C00: French Defence
w 39/11/0, b 4/46/0
3.♘c3 ♘f6 e5 ♘d7 f4 c5 ♘f3 ♗e7

B50: Sicilian Defence

B30: Sicilian Defence

# Deep networks with recurrent connections

*'a man sitting on a couch with a dog'*



Network desribes the image with the words:

*'a man sitting on a couch with a dog'*

(Fang et al. 2015)

# Quiz: Classification versus Reinforcement Learning

[ ] Classification aims at predicting the correct category such as 'car' or 'dog'

[ ] Classification is based on rewards

[ ] Reinforcement learning is based on rewards

[ ] Reinforcement learning aims at optimal action choices

[ ] Recurrent neural networks are useful for sequences

# Artificial Neural Networks: Lecture 1
## Simple Perceptrons for Classification

Wulfram Gerstner
EPFL, Lausanne, Switzerland

**Objectives for today:**
- understand classification as a geometrical problem
- discriminant function of classification
- linear versus nonlinear discriminant function
- perceptron algorithm
- gradient descent for simple perceptrons

# 1. The problem of Classification

car (yes or no)

output

the classifier

input

# 1. The problem of Classification

car (yes or no)

output

input

*Blackboard 1:*
 from images to vector

# 1. The problem of Classification

+1 yes (or 0 for no)

output ↑

the classifier
f($\mathbf{x}$)

input ↑ ↑ ↑ ↑ ↑ ↑

vector $\mathbf{x}$

# 1. Classification as a geometric problem



Blackboard 2:
from  vectors to classification

# 1. Classification as a geometric problem

**Task of Classification**

= find a **separating surface** in the high-dimensional input space

Classification by **discriminant function** $d(\boldsymbol{x})$

→ $d(\boldsymbol{x})=0$ on this surface; $d(\boldsymbol{x})>0$ for all positive examples $\boldsymbol{x}$

$d(\boldsymbol{x})<0$ for all counter examples $\boldsymbol{x}$

$d(\boldsymbol{x})=0$

$d(\boldsymbol{x}^{\mu}) > 0$

$d(\boldsymbol{x})=0$

linearly separable problem

# 2. Data base for Supervised learning

target output $t^\mu = 1$   **output**   $\hat{y}^\mu = 1$   classifier output

car (yes)

teacher

**Classifier**

input

$x^\mu$

## 2. Data base for Supervised learning

$P$ data points $\qquad \{ \quad (\boldsymbol{x}^{\mu}, t^{\mu}) \quad , \qquad 1 \leq \mu \leq P \quad \};$

input   target output

$$t^{\mu} = 1 \quad \text{car =yes}$$
$$t^{\mu} = 0 \quad \text{car =no}$$

# 2. Data base for Supervised learning

error!

car (no)

target output $t^7 = 1$ ⟷ $\hat{y}^7 = 0$   classifier output

**output**

**Classifier**

teacher

input

$x^7$

# 2. Data base for Supervised learning

$P$ data points $\quad \{ \quad (\boldsymbol{x}^{\mu}, t^{\mu}) \quad, \quad\quad 1 \leq \mu \leq P \quad \}$;

input  target output

for each data point $\boldsymbol{x}^{\mu}$, the classifier gives an output $\hat{y}^{\mu}$

→ **use errors** $\quad \hat{y}^{\mu} \neq t^{\mu}$ **for optimization of classifier**

**Remark**: for multi-class problems $\mathbf{y}$ and $\boldsymbol{t}$ are vectors

# 3. Single-Layer networks: simple perceptron

$$\hat{y} = g(a') = \begin{cases} +1 \text{ if } a' > \vartheta \\ 0 \text{ if } a' < \vartheta \end{cases}$$

$$\hat{y} = g\left(\sum_k w_k x_k\right)$$

$w_k$

$x_k$

output

$$\hat{y} = f(\boldsymbol{x})$$

the classifier
f($\boldsymbol{x}$)

vector $\boldsymbol{x}$

# 3. Single-Layer networks: simple perceptron

$$\hat{y}^{\mu} = 0.5[1 + sgn(\sum_k w_k\, x_k - \vartheta)]$$

output

$$\hat{y}^{\mu} = g\left(\sum_k w_k\, x_k\right)$$

$a$

$a'$

$w_{ik}$

$x_k$

$g(a')$

$\vartheta$

$a'$

$$g(a')= \begin{cases} 1 & if\ a' > \vartheta \\ 0.5 & if\ a' = \vartheta \\ 0 & if\ a' > \vartheta \end{cases}$$

input

vector $x$

# 3. Single-Layer networks: simple perceptron

$$\hat{y} = 0.5[1 + sgn(\sum_k w_k x_k - \vartheta)]$$

$$d(\boldsymbol{x}) = \sum_k w_k x_k - \vartheta = 0$$



$w_{ik}$

$x_k$

vector $\boldsymbol{x}$

imposes a linear separation

# 3. remove threshold: add a constant input

$$d(\boldsymbol{x}) = \sum_{k=1}^{N} w_k \, x_k - \vartheta = 0$$

$$d(\boldsymbol{x}) = \sum_{k=1}^{N+1} w_k \, x_k = 0$$

$w_{ik}$

$\boldsymbol{x} \in R^N$

$x_k$

$w_{N+1} = \vartheta$

$\boldsymbol{x} \in R^{N+1}$

$x_{N+1} = -1$

0

-1

# 3. Single-Layer networks: simple perceptron

## a simple perceptron

- can only solve linearly separable problems
- imposes a separating hyperplane
- for $\vartheta = 0$ hyperplane goes through origin
- threshold parameter $\vartheta$ can be removed by adding an input dimension
- in **N+1** dimensions hyperplane always goes through origin
- we can **adapt the weight vector** to the problem: this is called 'learning'

# 4. Perceptron algorithm: turn weight vector (in N+1 dim. )

$$hyperplane: d(\boldsymbol{x}) = \sum_{k=1}^{N+1} w_k\, x_k = \boldsymbol{w}^T \boldsymbol{x} = 0$$

# 4. Perceptron algorithm: turn weight vector

**Perceptron algo (in N+1 dimensions):**

   - set $\mu = 1$

  (1) cycle many times through patterns

   - choose pattern $\mu$

   - calculate output

$$\hat{y}^\mu = 0.5[1 + sgn(\boldsymbol{w}^T \boldsymbol{x}^\mu)]$$

   - update by

$$\Delta\boldsymbol{w} = \gamma[t^\mu - \hat{y}^\mu]\boldsymbol{x}^\mu$$

  - iterate $\mu \leftarrow (\mu + 1)mod\boldsymbol{P}$, back to (1)

  (2) stop if no changes for all $\boldsymbol{P}$ patterns

# 4. Perceptron algorithm: theoreom

If the problem is linearly separable, the perceptron algorithm converges in a finite number of steps.

Proof: in many books, e.g.,
Bishop, 1995,
*Neural Networks for Pattern Recognition*

# Quiz: Perceptron algorithm

The **input vector has *N* dimensions** and we apply a perceptron algorithm.
.

[ ] a rotation of the hyperplane in N+1 dimensions implies a change of weight vector.

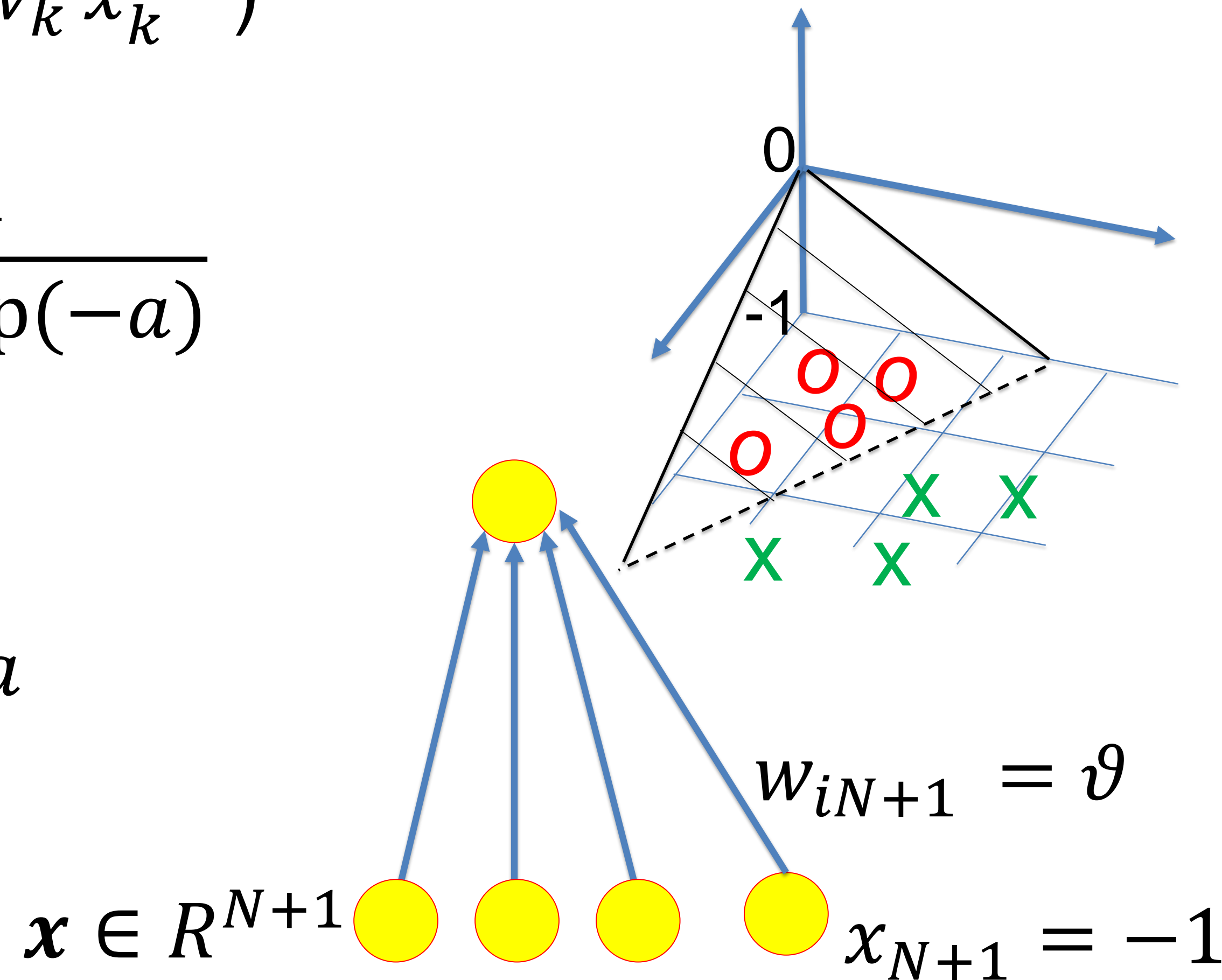[ ] An increase of the length of the weight vector implies that the hyperplane does not change in N+1 dimensions
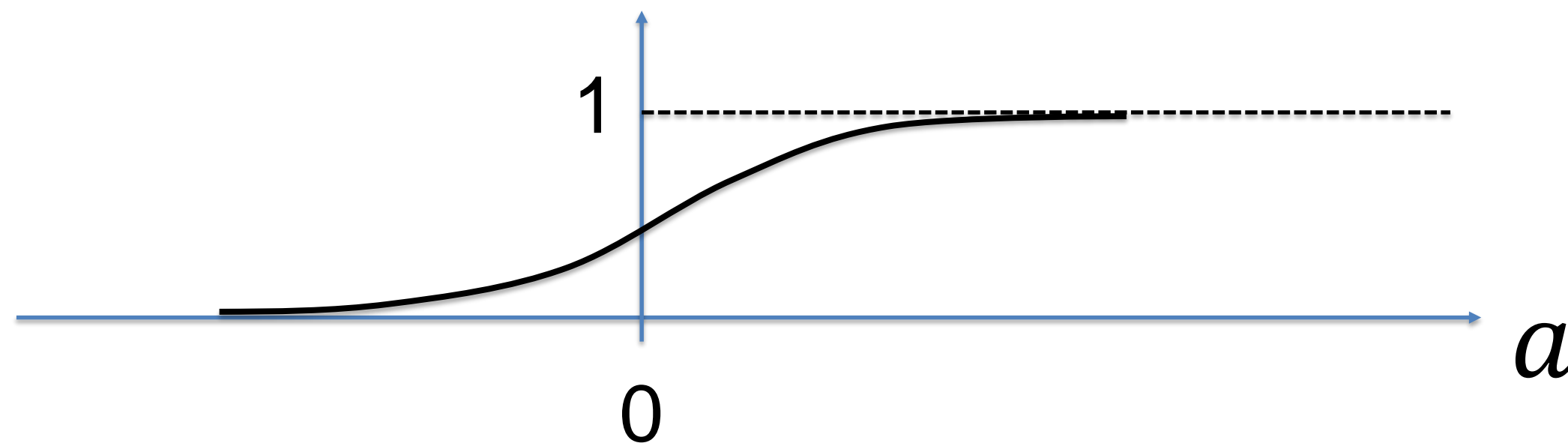
# 5. Sigmoidal output unit

A saturating nonlinear function with a smooth transition from 0 to 1.
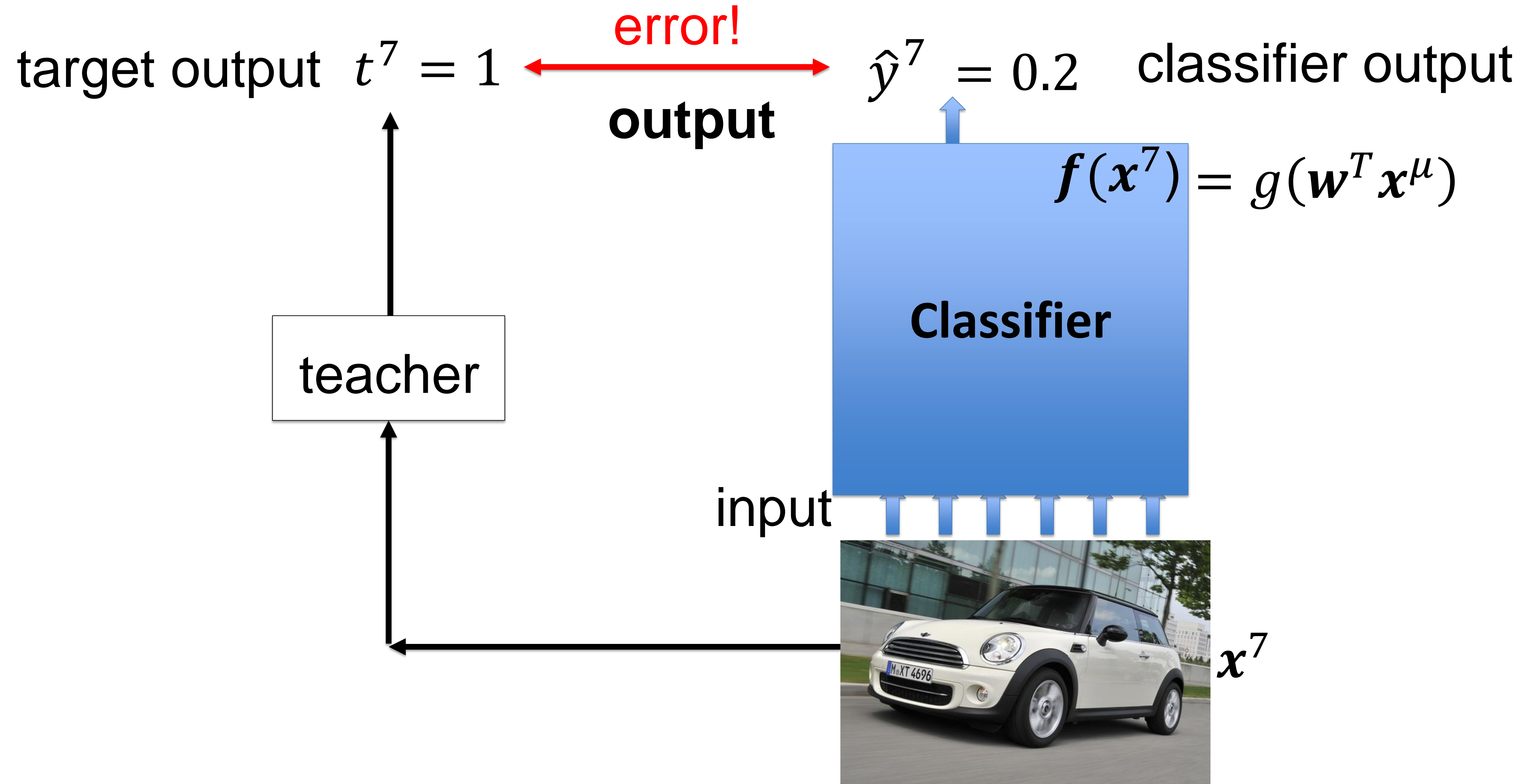
$$\hat{y}^{\mu} = g(\boldsymbol{w}^T \boldsymbol{x}^{\mu}) = g(\sum_{k=1}^{N+1} w_k \, x_k^{\mu})$$

with

$$g(a) = \frac{\exp(a)}{1 + \exp(a)} = \frac{1}{1 + \exp(-a)}$$

$\boldsymbol{x} \in R^{N+1}$
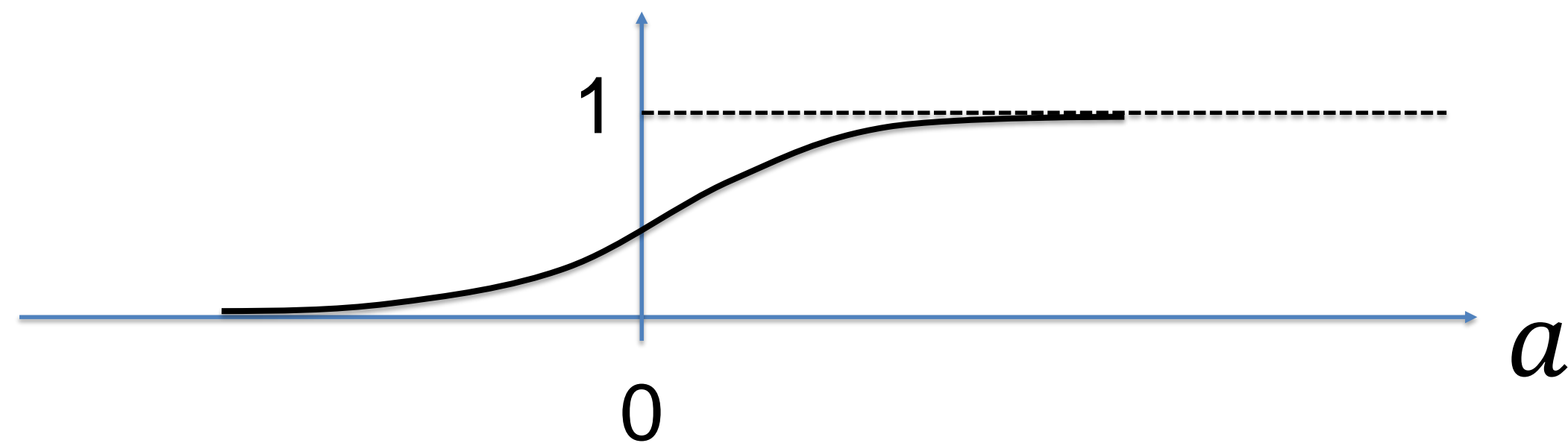
$w_{iN+1} = \vartheta$

$x_{N+1} = -1$

# 5. Supervised learning with sigmoidal output

error!

target output $t^7 = 1$ ⟷ $\hat{y}^7 = 0.2$   classifier output

**output**

$$f(x^7) = g(w^T x^\mu)$$

**Classifier**

teacher

input

$x^7$

# 5. Supervised learning with sigmoidal output

define **error**

$$E(\boldsymbol{w}) = \frac{1}{2}\sum_{\mu=1}^{P}\left[t^{\mu} - \hat{y}^{\mu}\right]^{2}$$



gradient descent

$$\Delta w_k = -\gamma \frac{dE}{dw_k}$$



$$\hat{y}^{\mu} = g(\boldsymbol{w}^{T}\boldsymbol{x}^{\mu})$$

$$w_{N+1} = \vartheta$$

$$\boldsymbol{x} \in R^{N+1}$$
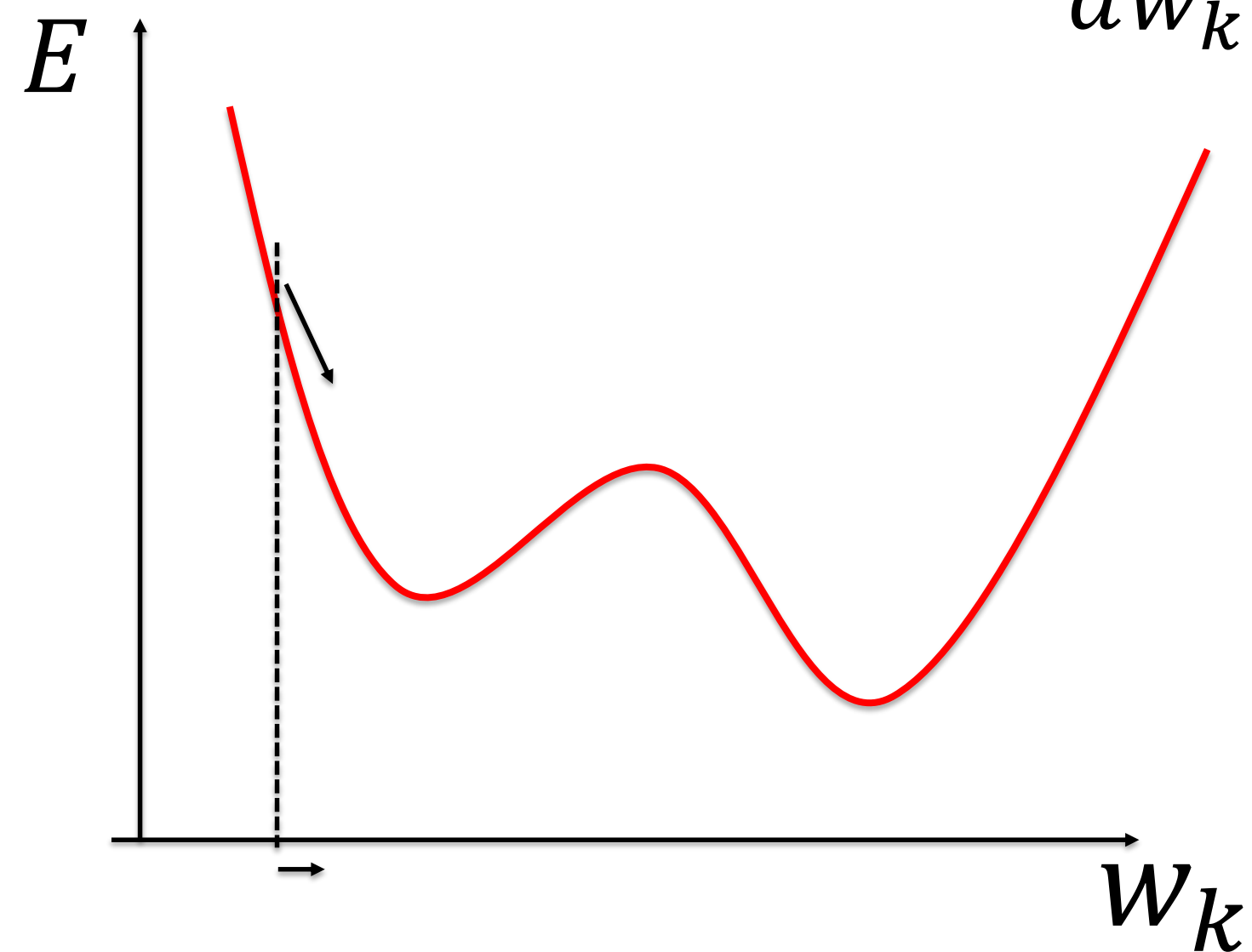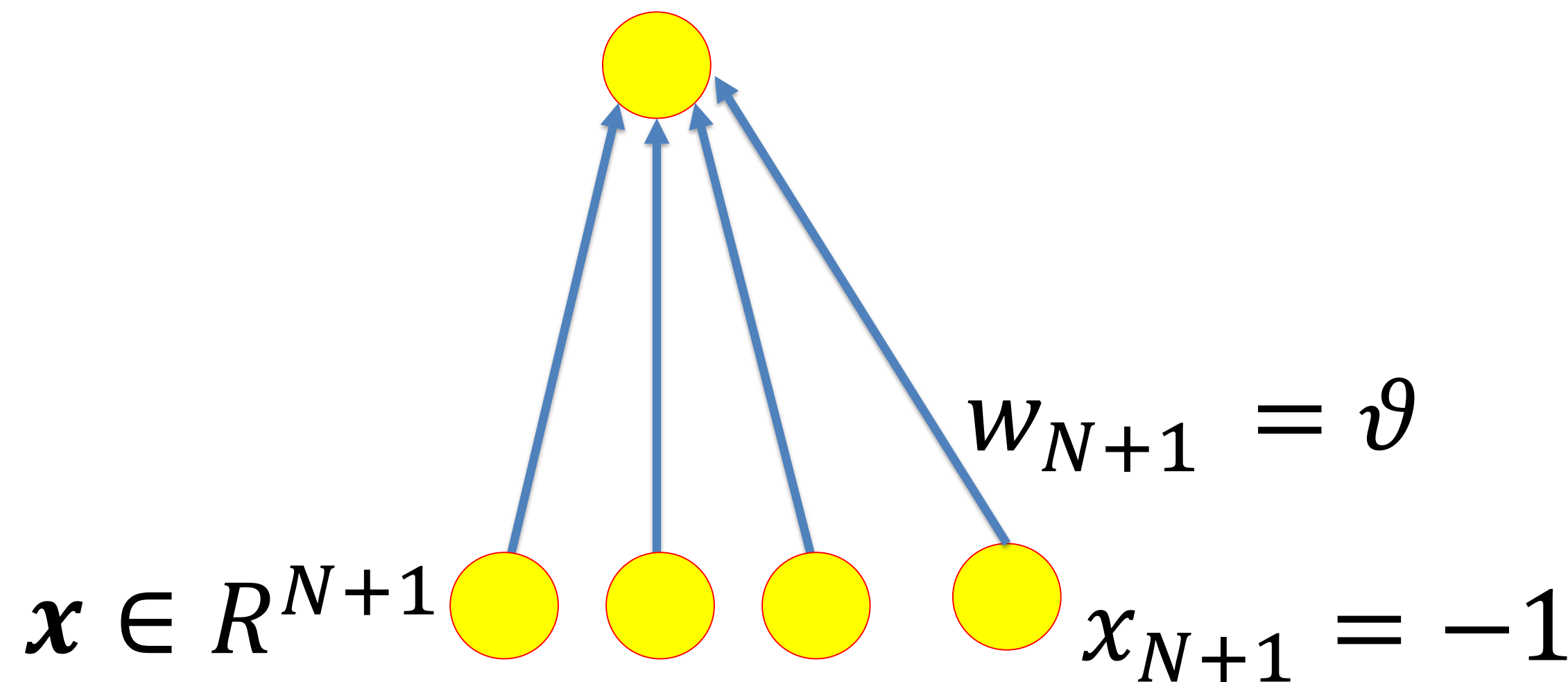
$$x_{N+1} = -1$$

# 6. gradient descent

Quadratic **error**

$$E(\boldsymbol{w}) = \frac{1}{2}\sum_{\mu=1}^{P}\left[t^{\mu} - \hat{y}^{\mu}\right]^2$$

gradient descent

$$w_k = -\gamma\frac{dE}{dw_k}$$



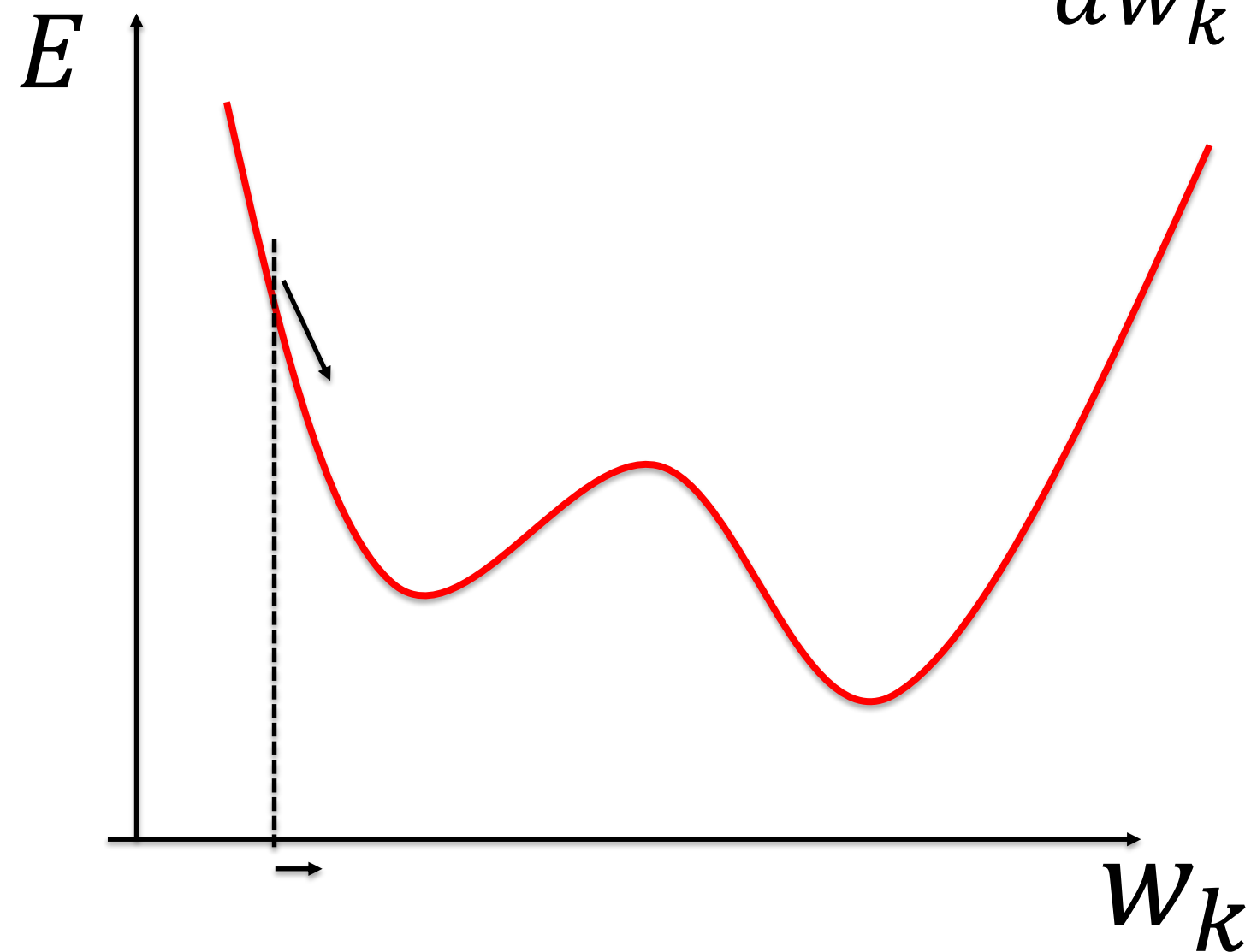$$\hat{y}^{\mu} = g(\boldsymbol{w}^T\boldsymbol{x}^{\mu})$$

$$w_{N+1} = \vartheta$$

$$\boldsymbol{x} \in R^{N+1}$$

$$x_{N+1} = -1$$

# 6. Gradient descent algorithm

After presentation of pattern $\boldsymbol{x}^\mu$ update the weight vector by

$$\Delta\boldsymbol{w} = \gamma\delta(\mu)\boldsymbol{x}^\mu$$

- amount of change depends on $\delta(\mu)$, i.e., the (signed) output mismatch for this data point
- change implemented even if 'correctly' classified
- change proportional to $\boldsymbol{x}^\mu$
- compare with perceptron algorithm

**Learning outcome and conclusions for today:**
- understand classification as a geometrical problem
- discriminant function of classification
- linear versus nonlinear discriminant function
- perceptron algorithm
- gradient descent for simple perceptrons
- learning as rotation of a hyperplane