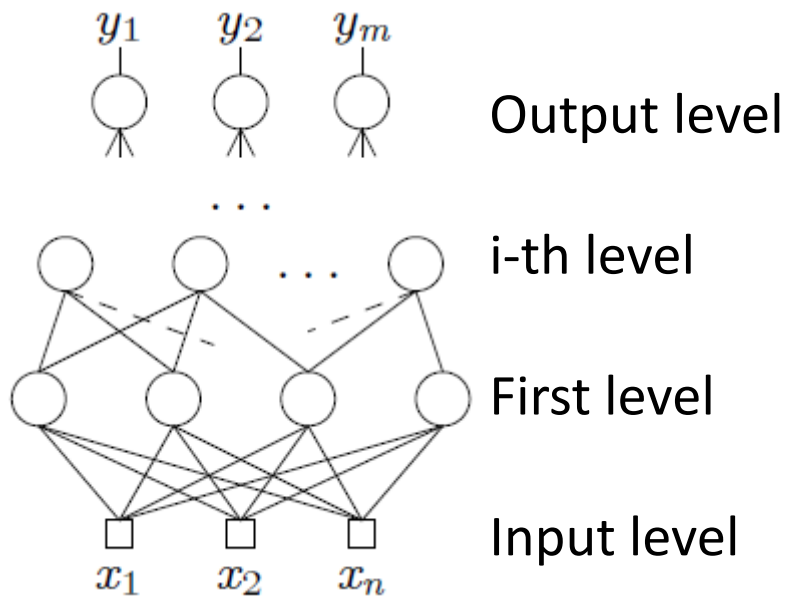


# Machine Learning

## Neural Networks

# Feed-Forward Networks



$i$  – Nummer der Schicht  
 $i = 0$  – Input Schicht  
 $i = i_{max}$  – Output Schicht  
 $j$  – Nummer des Neurons  
 $y_i$  – Output der  $i$ -ten Schicht (Vektor)  
 $y_{ij}$  – eine Komponente davon  
 $y_0 = x$  – Input Signal  
 $w_i$  – Gewichtsmatrix,  $i = 1 \dots i_{max}$   
 $w_{ijj'}$  – Gewicht, das  $(i, j)$  mit  $(i - 1, j')$  verbindet  
 $b_{ij}$  – Schwellwerte

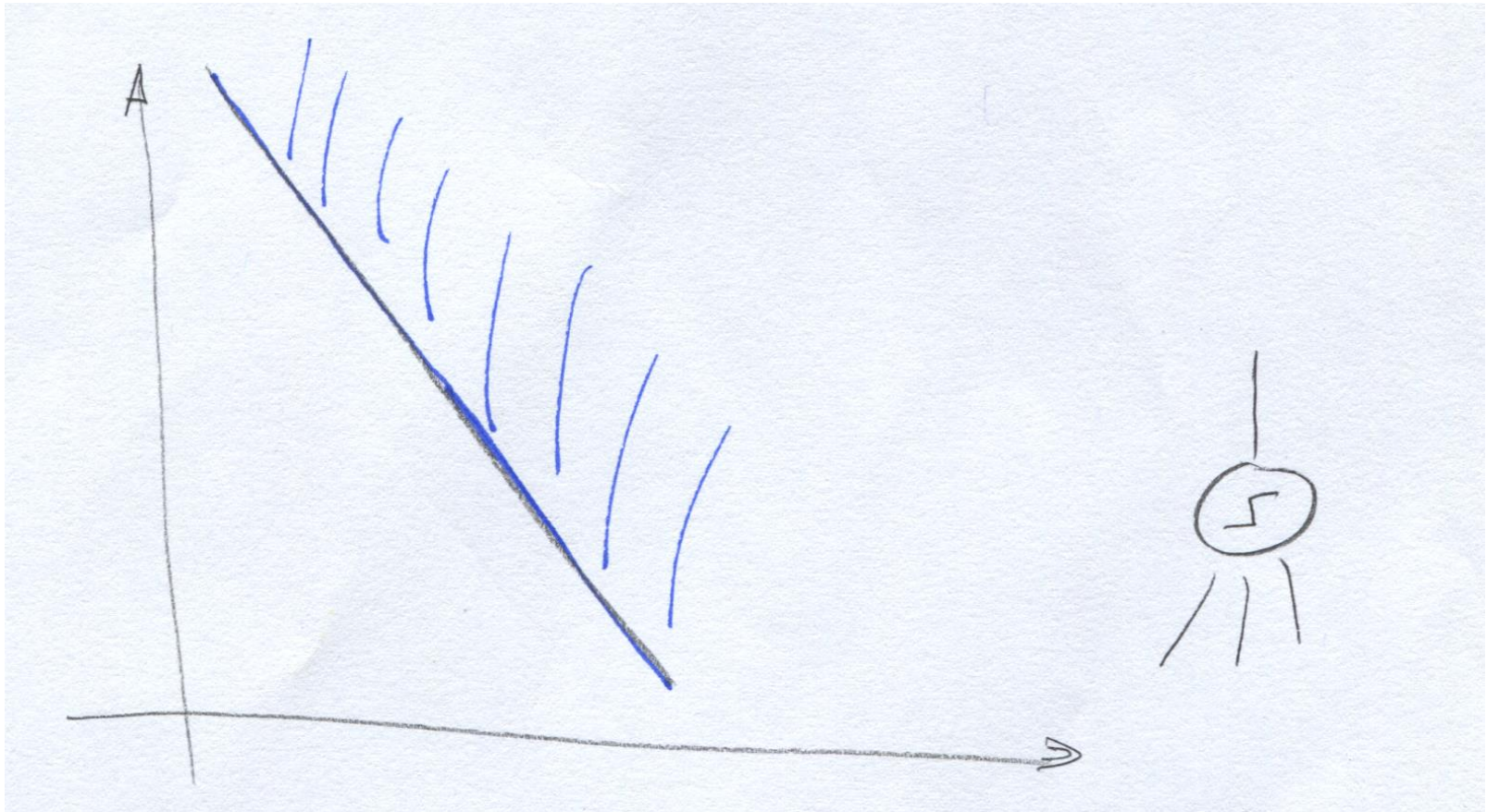
$$y_{ij} = f \left( \sum_{j'} w_{ijj'} y_{i-1 j'} - b_{ij} \right)$$

Special case:  $m = 1$  , Step-neurons – a mapping  $\mathbb{R}^n \rightarrow \{0, 1\}$

Which mappings can be modeled ?

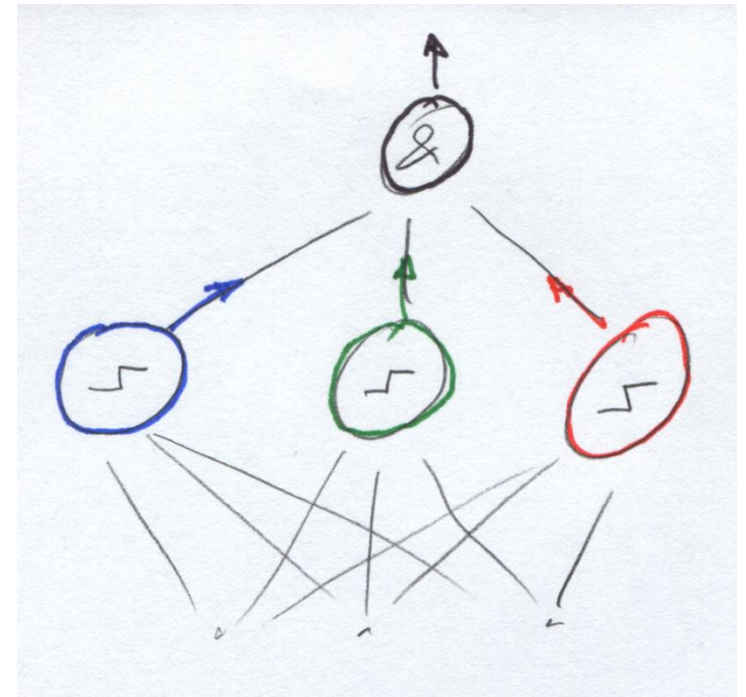
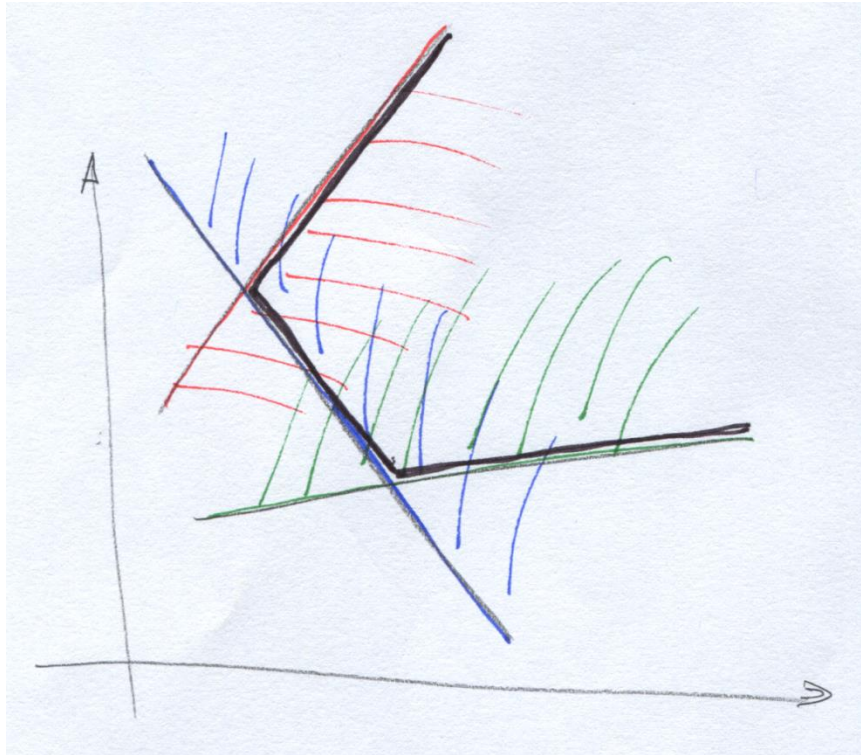
# Feed-Forward Networks

One level – single step-neuron – linear classifier



# Feed-Forward Networks

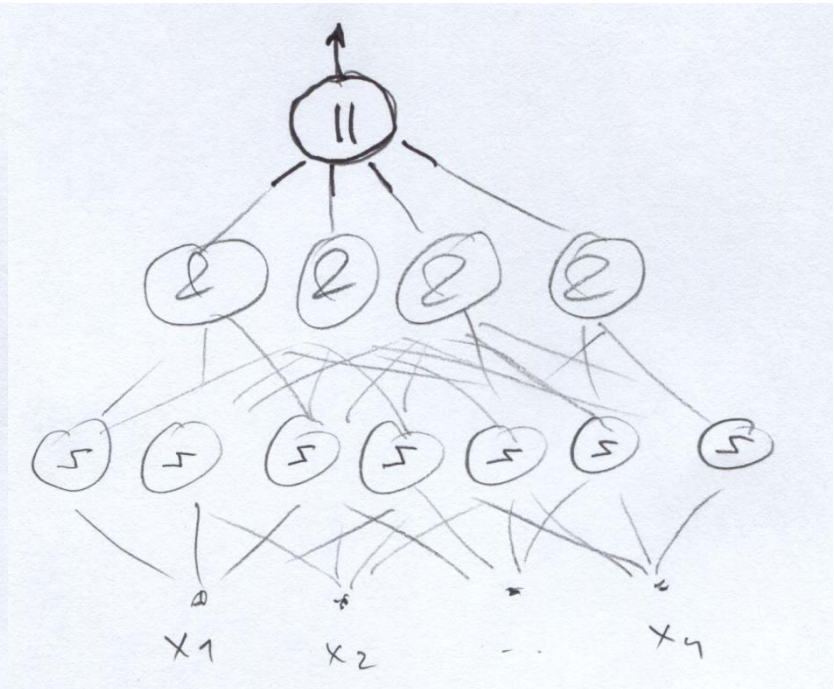
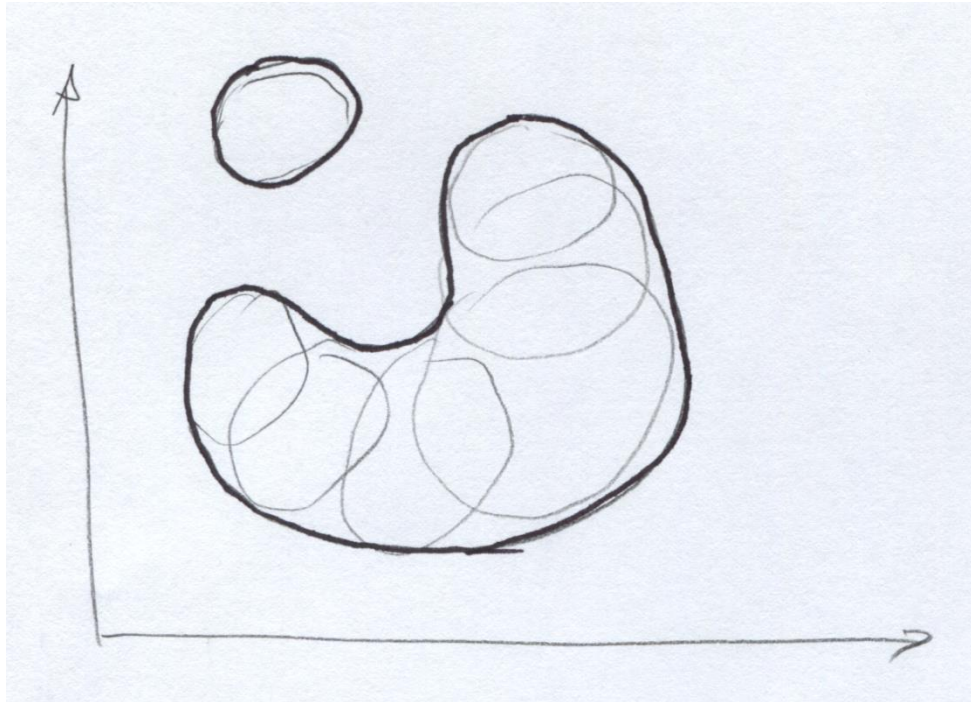
Two levels, “&”-neuron as the output – intersection of half-spaces



If the number of neurons is not limited, all convex subspaces can be implemented with an arbitrary precision.

# Feed-Forward Networks

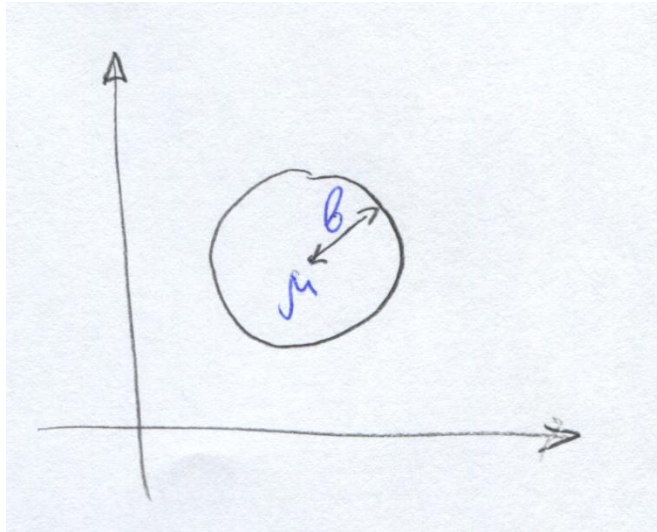
Three levels— all possible mappings  $\mathbb{R}^n \rightarrow \{0, 1\}$  as union of convex subspaces:



Three levels (sometimes even less) are enough to implement all possible mappings !!!

# Radial Basis Functions

Another type of neurons,  
corresponding classifier – “inside/outside a ball”



$$y = f(\|x - \mu\| - b)$$

The usage of RBF-neurons “replaces” a level in FF-networks.

With infinitely many RBF-neurons arbitrary mappings with only one intermediate level are possible.

# Error Back-Propagation

Learning task:

Given: training data  $((x^1, k^1) \dots (x^l, k^l))$ ,  $x^l \in \mathbb{R}^n$ ,  $k^l \in \mathbb{R}$

Find: all weights and biases of the net.

Error Back-Propagation is a **gradient descent** method for Feed-Forward-Networks with Sigmoid-neurons

First, we need an objective (error to be minimized)

$$F(w, b) = \sum_l (k^l - y(x^l; w, b))^2$$

Now: derive, build the gradient and go.

# Error Back-Propagation

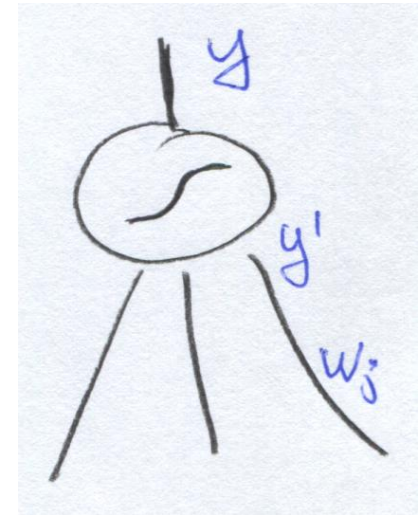
We start from a single neuron and just one example  $(x, k)$ .

Remember:

$$F(w, b) = (k - y)^2$$

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

$$y' = \langle x, w \rangle = \sum_j x_j w_j$$



Derivation according to the chain-rule:

$$\begin{aligned} \frac{\partial F(w, b)}{\partial w_j} &= \frac{\partial F}{\partial y} \cdot \frac{\partial y}{\partial y'} \cdot \frac{\partial y'}{\partial w_j} = \\ &= (y - k) \cdot \frac{\exp(-y')}{(1 + \exp(-y'))^2} \cdot x_j = \delta \cdot d(y') \cdot x_j \end{aligned}$$



# Error Back-Propagation

The “problem”: for intermediate neurons the errors are not known !  
 Now a bit more complex:

$$\frac{\partial F}{\partial w_*} = \frac{\partial F}{\partial y_n} \cdot \frac{\partial y_n}{\partial y'_n} \cdot \left[ \sum_j \frac{\partial y'_n}{\partial y_j} \cdot \frac{\partial y_j}{\partial y'_j} \cdot \frac{\partial y'_j}{\partial y_*} \right] \cdot \frac{\partial y_*}{\partial y'_*} \cdot \frac{\partial y'_*}{\partial w_*} =$$

$$\left[ \sum_j \delta_n \cdot d(y'_n) \cdot w_{nj} \cdot d(y'_j) \cdot w_{j*} \right] \cdot d(y'_*) \cdot x =$$

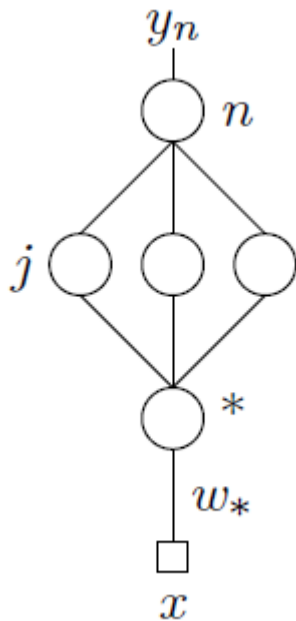
$$\left[ \sum_j \delta_j \cdot d(y'_j) \cdot w_{j*} \right] \cdot d(y'_*) \cdot x =$$

$$\delta_* \cdot d(y'_*) \cdot x$$

$$\delta_j = \delta_n \cdot d(y'_n) \cdot w_{nj}$$

$$\delta_* = \left[ \sum_j \delta_j \cdot d(y'_j) \cdot w_{j*} \right]$$

with:



# Error Back-Propagation

In general: compute “errors”  $\delta$  at the  $i$ -th level from all  $\delta$ -s at the  $i+1$ -th level – propagate the error.

The Algorithm (for just one example  $(x, k)$ ):

1. Forward: compute all  $y$  and  $y'$  (apply the network), compute the output error  $\delta_n = y_n - k$ ;
2. Backward: compute errors in the intermediate levels:

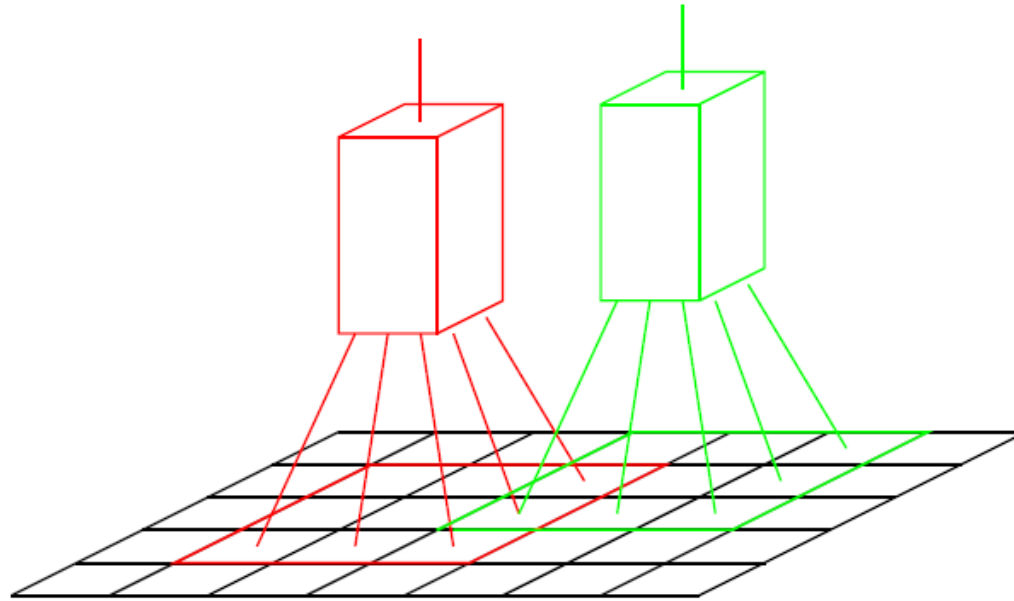
$$\delta_{ij} = \sum_{j'} \delta_{i+1 j'} \cdot d(y_{i+1 j'}) \cdot w_{i+1 j' j}$$

3. Compute the gradient and go.

$$\frac{\partial F}{\partial w_{ijj'}} = \delta_{ij} \cdot d(y'_{ij}) \cdot y_{i-1 j'}$$

For many examples – just sum them up.

# Time Delay Neural Networks (TDNN)



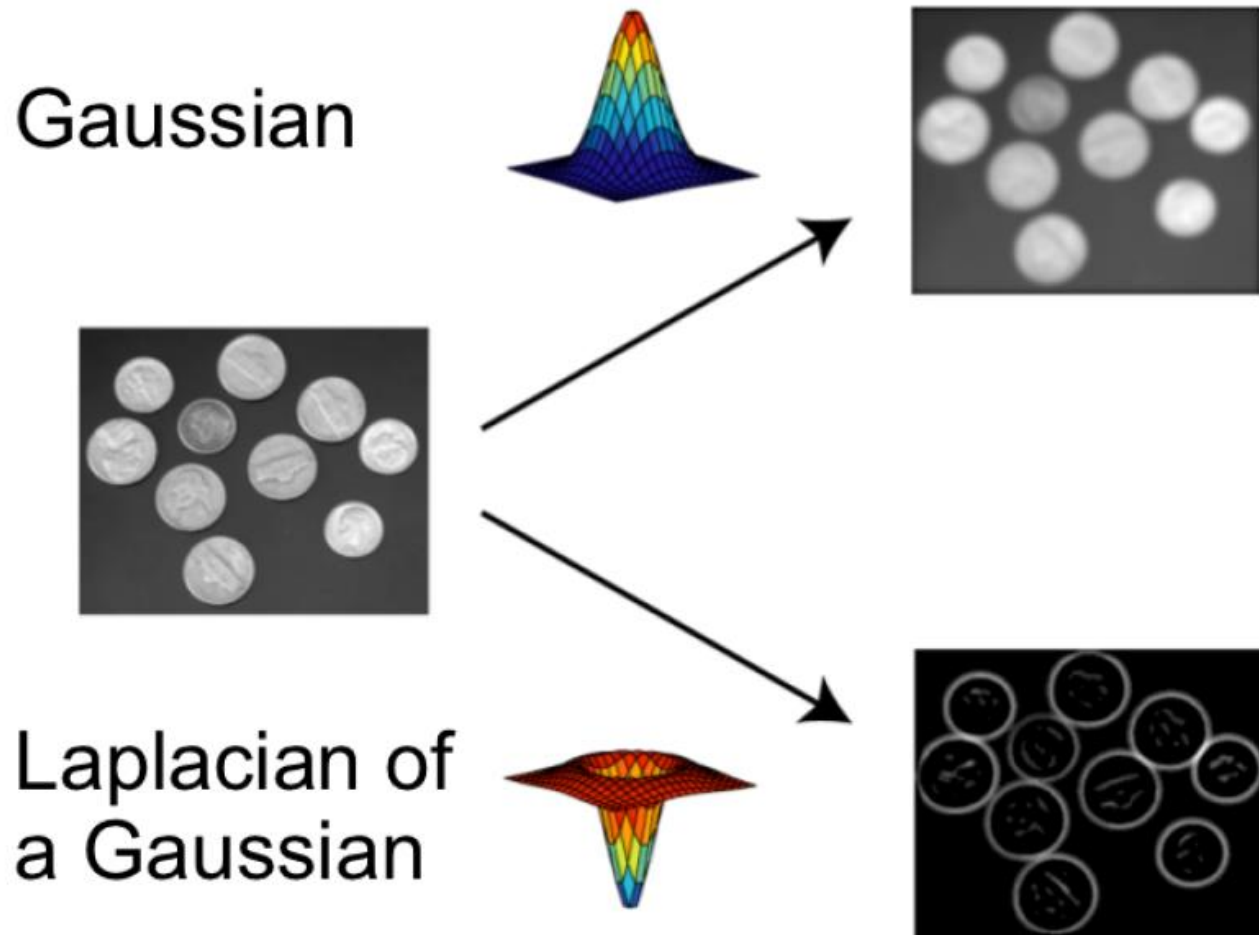
Feed-Forward network of a particular architecture.

Many equivalent “parts” (i.e. of the same structure with the same weights), but having different **Receptive Fields**. The output level of each part gives an information about the signal in the corresponding receptive field – computation of **local features**.

Problem: During the Error Back-Propagation the equivalence gets lost. Solution: average the gradients.

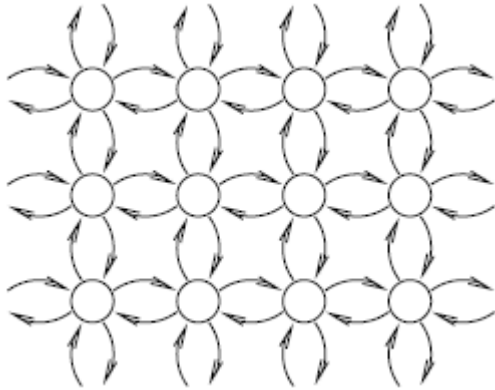
# Convolutional Networks

Local features – convolutions with a set of predefined masks (see lectures “Computer Vision”).



# Hopfield Networks

There is a symmetric neighborhood relation (e.g. a grid).  
The output of each neuron serves as inputs for the neighboring ones.



$$y_r = f \left( \sum_{r' \in N(r)} w_{rr'} \cdot y_{r'} + b_r \right)$$

with symmetric weights, i.e.  $w_{rr'} = w_{r'r}$

A **network configuration** is a mapping  $y : D \rightarrow \{0, 1\}$

A configuration is **stable** if “outputs  $y_r$  do not contradict”

The **Energy** of a configuration is

$$E(y) = \sum_{rr'} w_{rr'} \cdot y_r \cdot y_{r'} + \sum_r b_r \cdot y_r$$

# Hopfield Networks

Network dynamic:

1. Start with an arbitrary configuration  $y^{(0)}$ ,
2. Decide for each neuron whether it should be activated or not according to

$$y_r = f \left( \sum_{r' \in N(r)} w_{rr'} \cdot y_{r'} + b_r \right)$$

Do it **sequentially** for all neurons until convergence, i.e. apply the changes immediately.

In doing so the energy increases !!!

Attention!!! It does not work with the parallel dynamic (seminar).

# Hopfield Networks

During the sequential dynamic the energy may only increase !

Proof:

Consider the energy “part” that depend on a particular neuron:

$$E_r(y) = \sum_{r' \in N(r)} w_{rr'} \cdot y_r \cdot y_{r'} + b_r \cdot y_r = y_r \cdot \left[ \sum_{r' \in N(r)} w_{rr'} \cdot y_{r'} + b_r \right]$$

After the decision the energy difference is

$$\begin{aligned} E^{(t+1)}(y) - E^{(t)}(y) &= \\ &= \left( y_r^{(t+1)} - y_r^{(t)} \right) \cdot \left[ \sum_{r' \in N(r)} w_{rr'} \cdot y_{r'} + b_r \right] \end{aligned}$$

If  $[\cdot] > 0$  , the new output  $y_r^{(t+1)}$  is set to 1  $\rightarrow$  energy grows.

If  $[\cdot] < 0$  , the new output  $y_r^{(t+1)}$  is set to 0  $\rightarrow$  energy grows too.

# Hopfield Networks

The network dynamic is the simplest method to find a configuration of the maximal energy (synonym – “Iterated Conditional Modes”).

The network dynamic is **not globally optimal**, it stops at a **stable** configuration, i.e. a **local** maxima of the Energy.

The **most stable** configuration – **global** maximum.

The task (find the global maximum) is NP-complete in general.

Polynomial solvable special cases:

1. The neighborhood structure is simple – e.g. a tree
2. All weights  $w_{TT'}$  are non-negative (supermodular energies).

Of course, nowadays there are many good approximations.



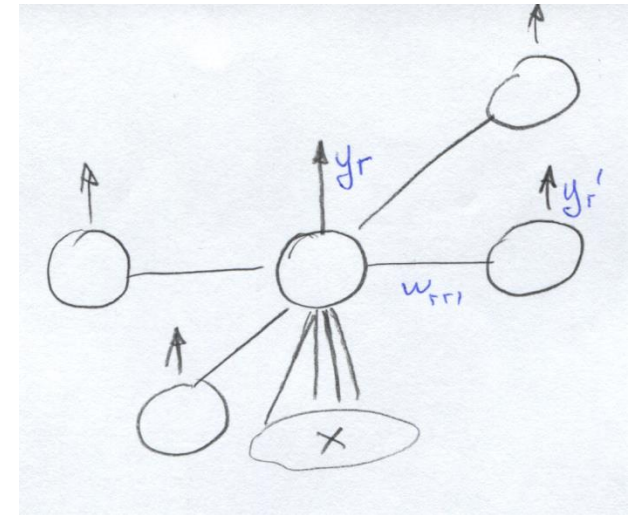
# Hopfield Networks

Hopfield Network with external input  $x$ :

$$y_r = f\left(\sum_{r' \in N(r)} w_{rr'} \cdot y_{r'} + b_r(x)\right)$$

The energy is

$$E(y, x) = \sum_{rr'} w_{rr'} \cdot y_r \cdot y_{r'} + \sum_r b_r(x) \cdot y_r$$



Hopfield Networks implement mappings  $X \rightarrow Y$  according to the principle of Energy maximum:

$$y = \arg \max_{y'} E(y', x)$$

Note: no single output but a configuration – **structured output**.

# Hopfield Networks

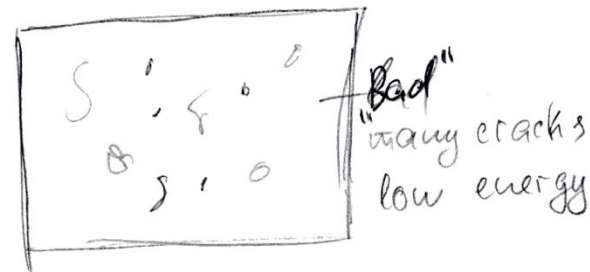
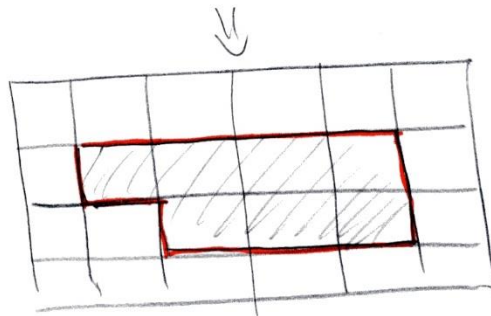
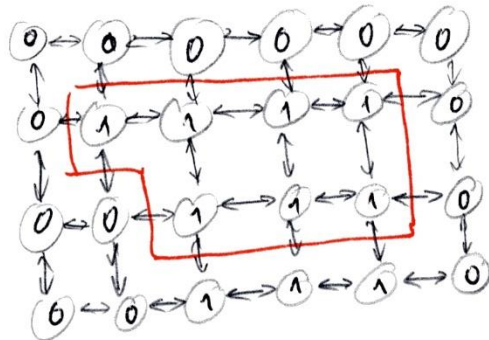
Hopfield Networks model patterns  
– network configurations of the optimal energy.

Example:

Let  $y$  be a network configuration and  $n(y)$  the number of “cracks” – pairs of neighboring neurons of different outputs.

Design a network (weights and biases for each neuron) so that the energy of a configuration is proportional to the number of cracks, i.e.  $E(y) \sim -n(y)$  .

# Hopfield Networks



Solution:  $w_{rr'} = 2$ ,  $b_r = -4$  (up to the borders)

Further examples at the seminar.