# ARTIFICIAL NEURAL NETWORKS

Federico Marini
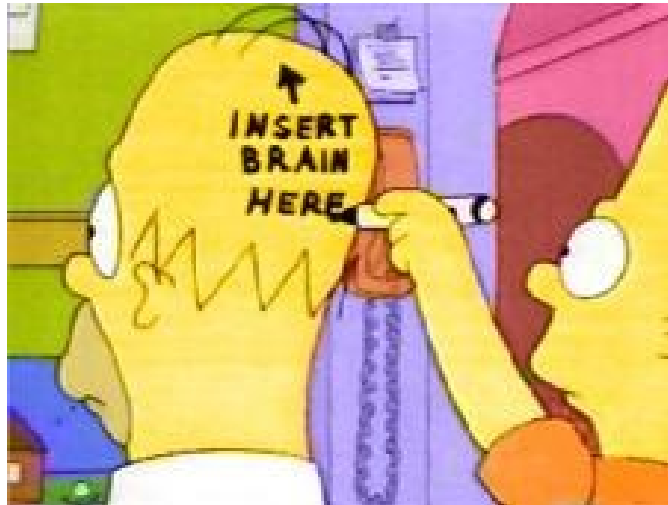
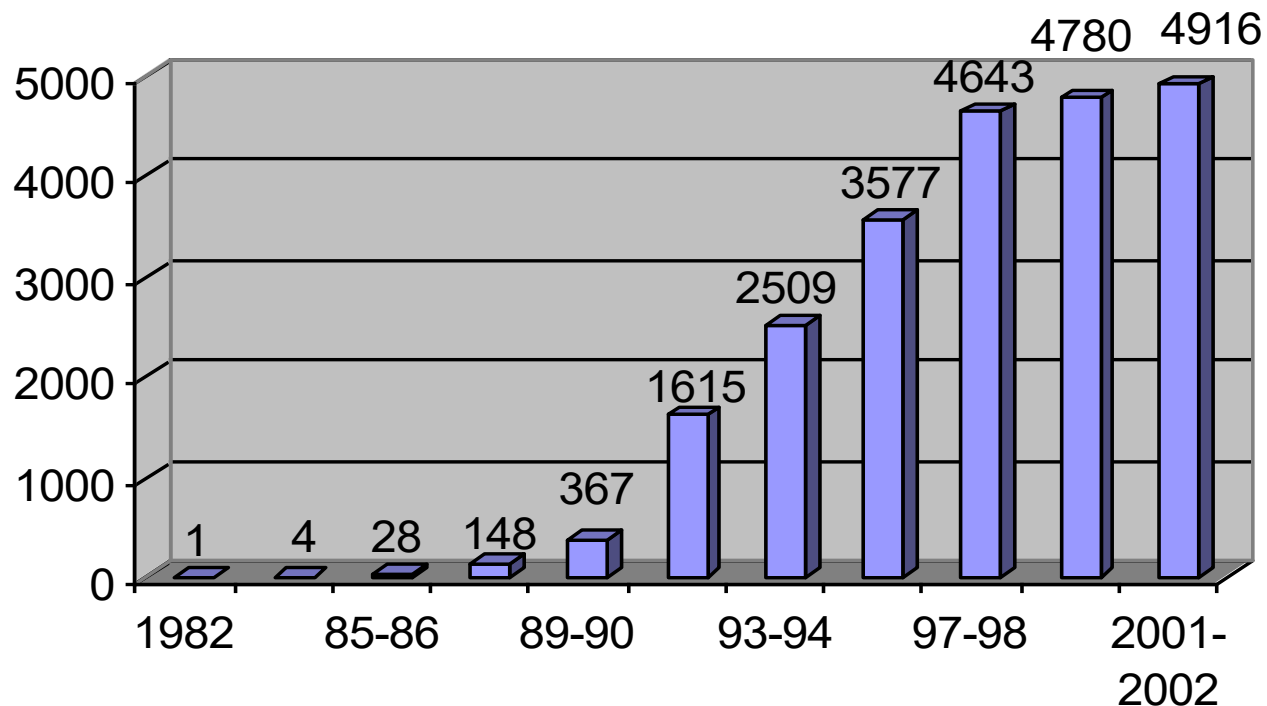*Dept. Chemistry, University of Rome "La Sapienza", Italy*

# Artificial neural networks (ANNs)

# Historical background

- 1943: McCulloch & Pitts model of a neuron
- 1949: Hebb postulates a learning rule
- 1958: Rosenblatt model of a neural net
- 1969: the book "Perceptrons" freezes most of the initial enthusiasm about ANNs
- 1982: Hopfield and its model of neural net act as "catalyst" in attracting the attention of many scientists towards NNs
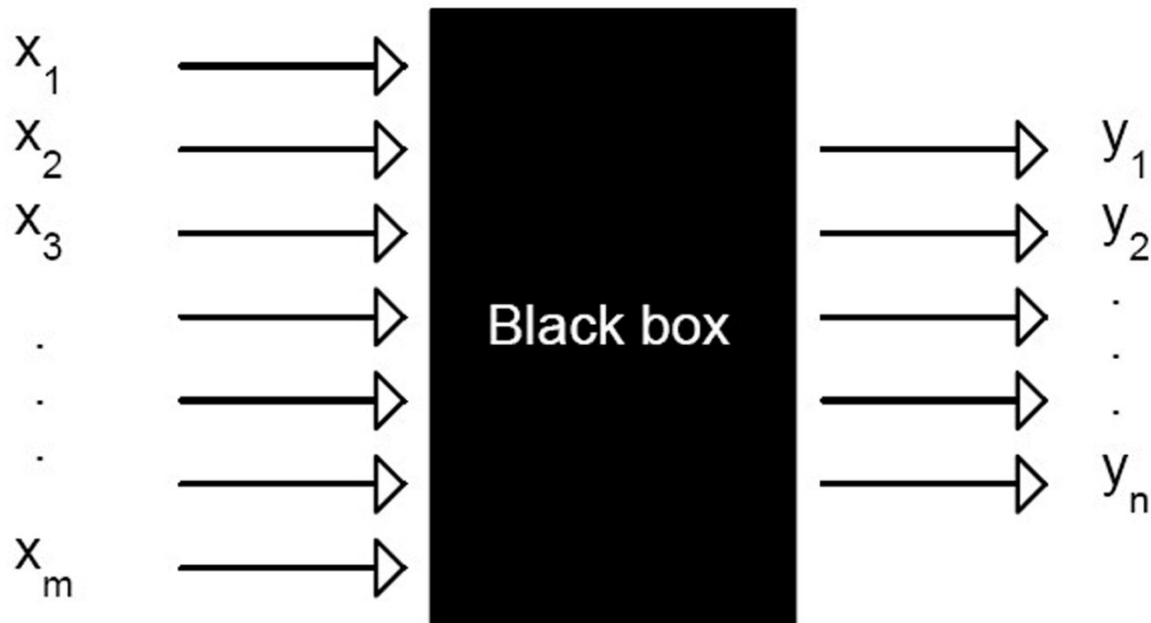
ANN papers published: 1982-2002

J. Hopfield, Neural Networks and Physical Systems with Emergent Collective Computational Ability, P. Nac. Acad. Sci. Biol., 79(2), (1982), 2554-2558

# Strengths of a Neural Network

- **Power**: Model complex functions, nonlinearity built into the network
- **Ease of use**:
  - Learn by example
  - Very little user domain-specific expertise needed
- **Intuitively appealing**: based on model of biology, will it lead to genuinely intelligent computers/robots?

Neural networks cannot do anything that cannot be done using traditional computing techniques, **BUT** they can do some things which would otherwise be very difficult.

# The roughest approach to NNs



$x_1$
$x_2$
$x_3$
.
.
.
$x_m$

Black box

$y_1$
$y_2$
.
.
.
$y_n$

Input variables          Non-linear relation          Output varibles

# NNs in a nutshell

- From a computational point of view, ANNs represent a way to operate a non-linear functional mapping between an input and an output space.

$$\mathbf{y} = f(\mathbf{x})$$

- **Y** can be:
  - an n-D (usually 2D) vector of coordinates (*mapping*)
  - A multi-dimensional vector of response (*regression*)
  - A binary vector of class-memberships (*classification*)
- This functional relation is expressed in an <u>implicit</u> way

# **Opening the blackbox**

- The peculiarity of ANNs relies on the fact that they operate using a large number of parallel connected simple arithmetic units (neurons).

- Mathematically, a neuron can be defined as *nonlinear, parameterized, bounded function* ➔ the variables this function depends on are called the inputs of the neuron and its value is called the output.

$$y = f\left(\sum_i w_i x_i + w_0\right)$$

# Artificial neurons

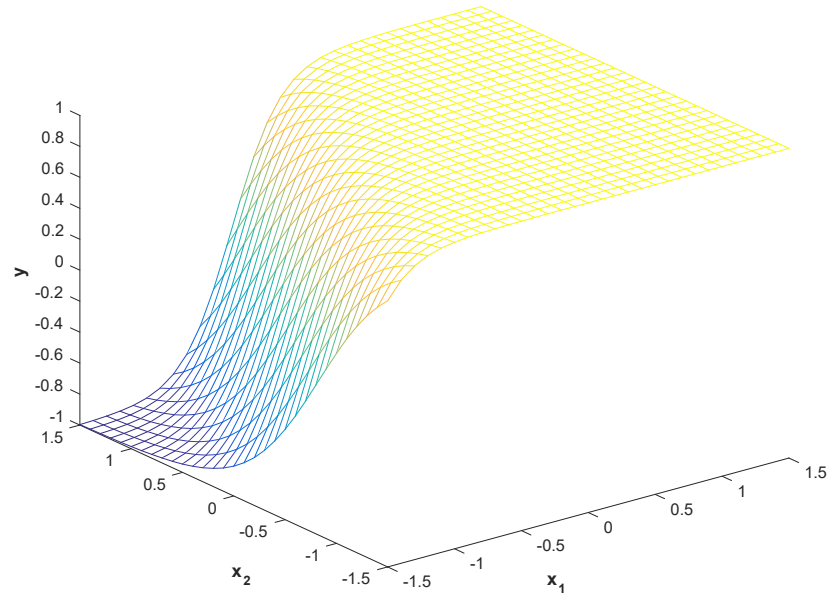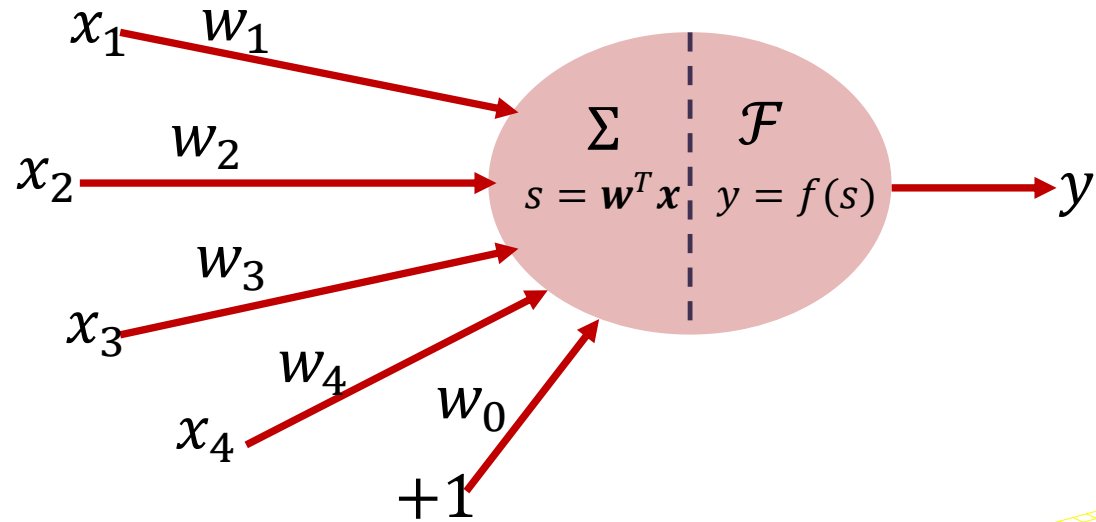Nonlinear generalization of the McCullogh-Pitts neuron:

$$y = f(x, w)$$

y is the neuron's output, x is the vector of inputs, and w is the vector of synaptic weights.

Examples:

$$y = \frac{1}{1 + e^{-w^T x - a}}$$ sigmoidal neuron

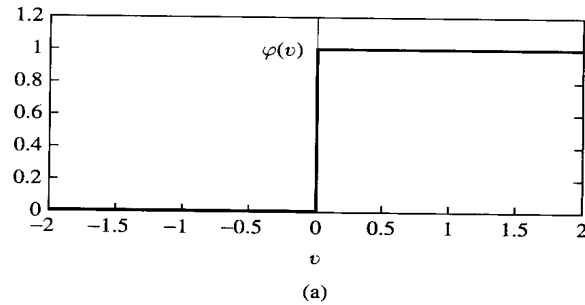$$y = e^{-\frac{||x-w||^2}{2a^2}}$$ Gaussian neuron

# Artificial neurons $y = f(x, w)$

$x_1$ $w_1$

$x_2$ $w_2$

$w_3$

$x_3$

$w_4$

$x_4$

$w_0$

$+1$

$\Sigma$  $\mathcal{F}$

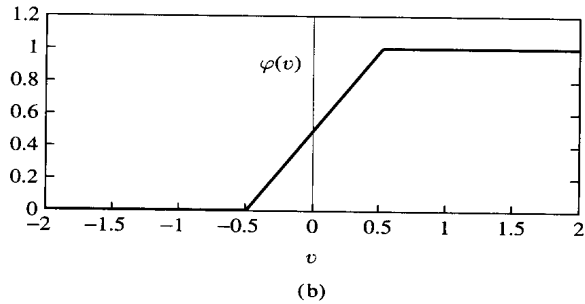$s = \boldsymbol{w}^T \boldsymbol{x}$ | $y = f(s)$

$y$

# Activation functions

Hard threshold
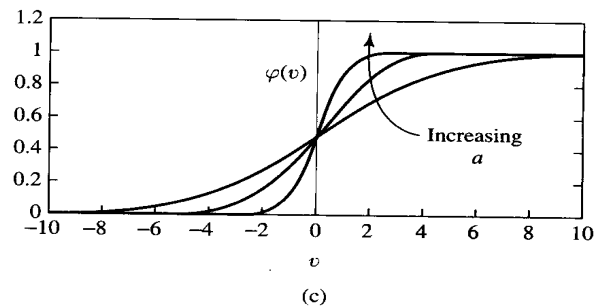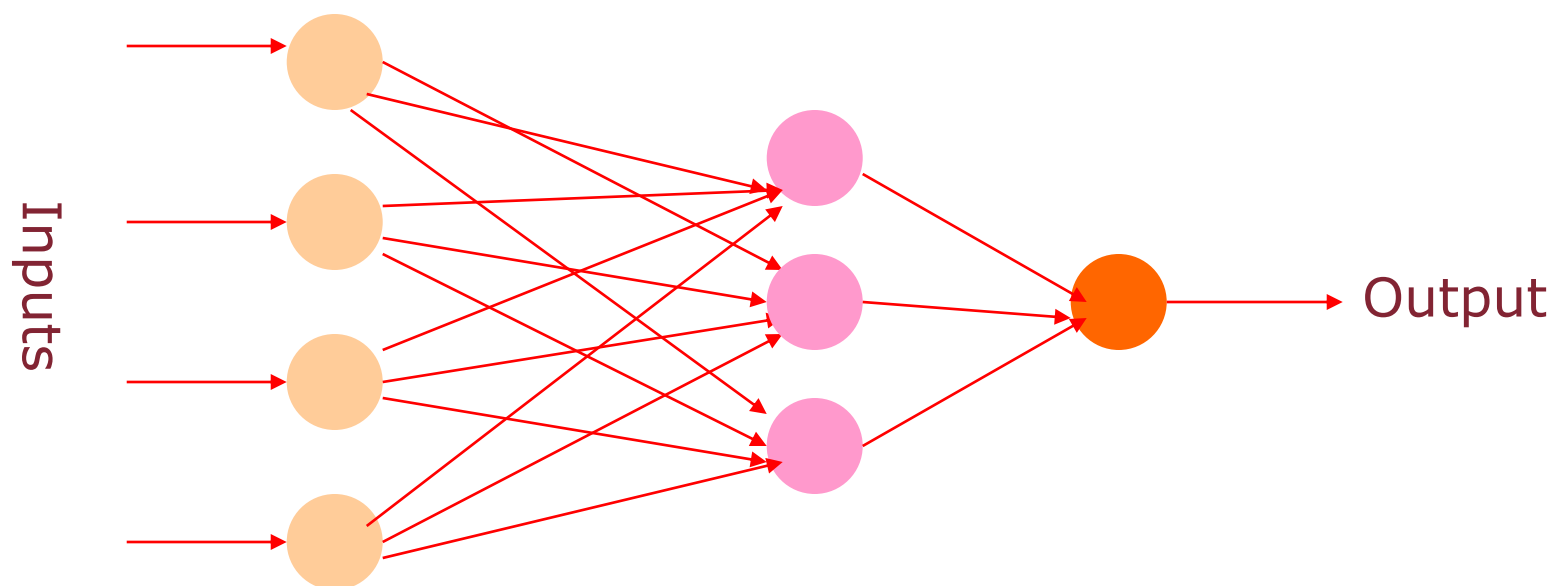
Piecewise linear

Sigmoid



**FIGURE 1.8** (a) Threshold function. (b) Piecewise-linear function. (c) Sigmoid function for varying slope parameter *a*.
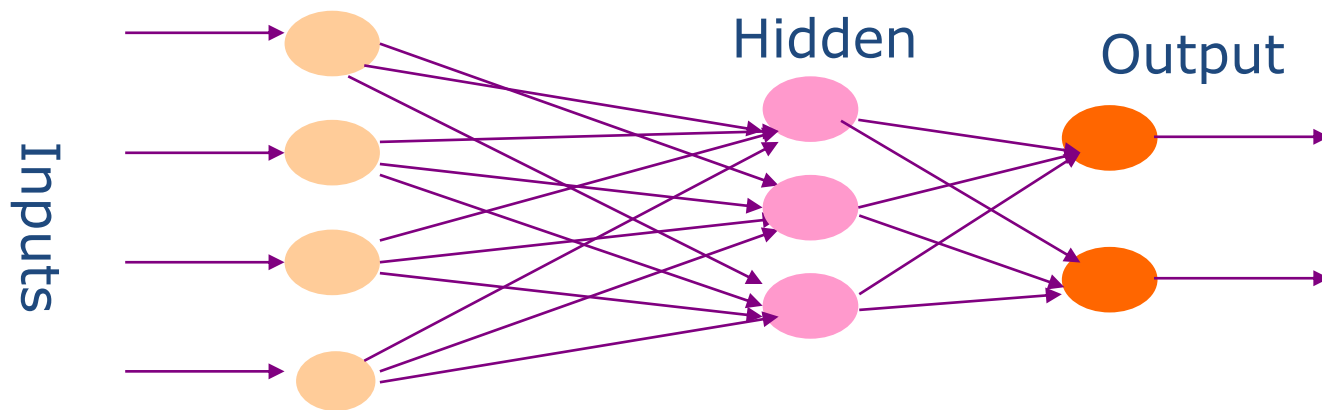
# From the neuron to the net

- Just as a neuron can be thought as a nonlinear function of its inputs, a network represents the composition of the nonlinear functions of two or more neurons

- The way the different units are connected among each other governs the way the different functions they describe are weighted and combined to produce the overall output.

- This pattern of interconnection among the neurons is called the network **"architecture"**, and can be conveniently represented on a graph: neurons operating on the same input variables are organized in *layers*, while the weights that modulate the combination of the nonlinear functions are represented as lines connecting units in different layers.

# Artificial neural networks

Inputs

Output

An artificial neural network is composed of many artificial neurons that are linked together according to a specific network architecture. The objective of the neural network is to transform the inputs into meaningful outputs.

# Artificial (multilayer feed-forward) NNs



Inputs

Hidden

Output

$$y_j = g\left(\sum_{k=1}^{N_h} w_{jk} h_k + w_{j0}\right) = g\left(\sum_{k=1}^{N_h} w_{jk} f\left(\sum_{i=1}^{N_i} w_{ki} x_i + w_{k0}\right) + w_{j0}\right)$$

# Artificial neural networks



$s = \text{sum}(w^T x)$

$h = \tanh(s)$

$y = \tanh(w_2^T h)$

# How does a neural network learn?

- A neural network learns by determining the relation between the inputs and outputs.

- By calculating the relative importance of the inputs and outputs the system can determine such relationships.

- Through trial and error the system compares its results with the expert provided results in the data until it has reached an accuracy level defined by the user.
  - With each trial the weight assigned to the inputs is changed until the desired results are reached.

# Multilayer feed-forward networks

- The most common network architecture is the **feed-forward** one, in which the information flows only in the forward direction, from inputs to outputs.

- This means that its graph representation is acyclic: no path in the graph, following the connections, can lead back to the starting point.

- A great variety of network topologies can be imagined, under the sole constraint that the graph of connections be acyclic. However, as anticipated before, the vast majority of neural network applications implement multilayer networks.

# Multilayer feed-forward NNs 2

Inputs

Hidden

Output

The network computes as many functions of the input variables of the network as are the components of the output vector

• Neurons are organized in three kind of layers: input, hidden and output.

• The output neurons are the neurons that perform the final computation, i.e., whose outputs are the outputs of the network, while the other neurons, which perform intermediate computations, are termed hidden neurons.

• The units of the input layer just pass the inputs as variables to the hidden neurons, without doing any processing on them.

• Each output is a nonlinear function (computed by the corresponding output neuron) of the nonlinear functions computed by the hidden neurons.

$$y_j = g\left(\sum\nolimits_{k=1}^{N_h} w_{jk} h_k + w_{j0}\right) = g\left(\sum\nolimits_{k=1}^{N_h} w_{jk} f\left(\sum\nolimits_{i=1}^{N_i} w_{ki} x_i + w_{k0}\right) + w_{j0}\right)$$
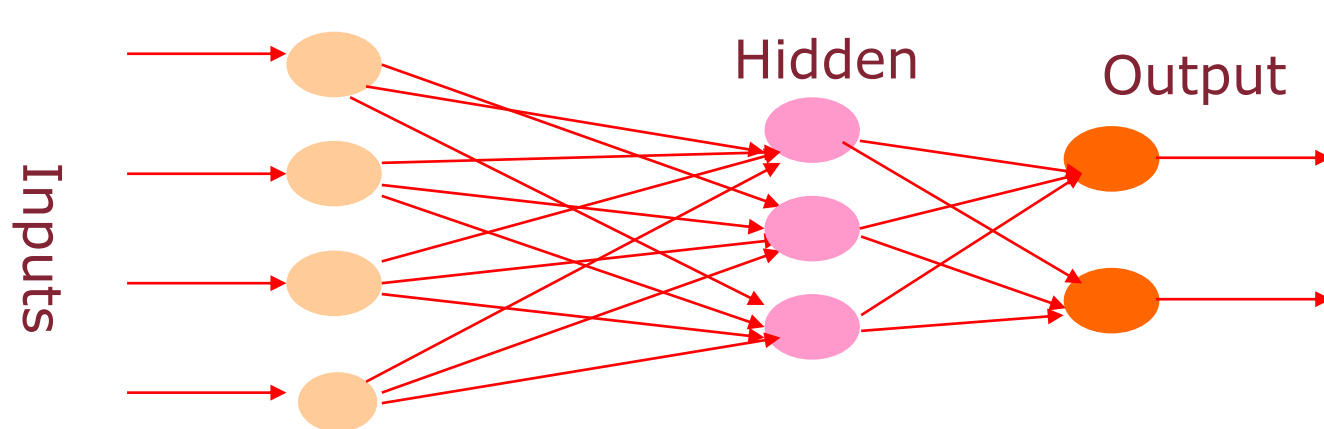
# "Training" the net

- A feed-forward network with a single hidden layer can approximate with arbitrary accuracy any bounded and sufficiently regular function in a finite region of variable space.

- The procedure by whereby the parameters of the network are estimated, in order to approximate such a function is called *training* of the ANN.

- Since usually the nonlinear relationship between dependent and independent variables is not known analytically, but a finite number of numerical values of the function are known (because they are obtained through measurements performed on a physical, chemical, biological, etc. process): *the task that is assigned to the network is that of approximating the regression function of the available data*.

# Supervised training

- This kind of training is referred to as "supervised" since the function that the network should implement is known in some or all points: a "teacher" provides "examples" of values of the inputs and of the corresponding values of the output,

- The goal of the training algorithm is to find the best set of model parameters, given the data ➔ find the numerical values of the network weights which minimize a cost function representing the distance between the prediction of the model and the measured values.

$$E^P(\mathbf{w}) = \frac{1}{2}\sum_{k=1}^{N_o}(y_k - \hat{y}_k)^2 = \frac{1}{2}\|\mathbf{y} - \hat{\mathbf{y}}\|^2$$

# Supervised training

- When this cost function is the squared error of the residuals:

$$E^{\text{batch}}(\mathbf{w}) = \frac{1}{N_{\text{samples}}} \sum_{P=1}^{N_{\text{samples}}} E^P(\mathbf{w}) = \frac{1}{2N_{\text{samples}}} \sum_{P=1}^{N_{\text{samples}}} \sum_{k=1}^{N_o} \left(y_{k,P} - \hat{y}_{k,P}\right)^2 = \frac{1}{2N_{\text{samples}}} \sum_{P=1}^{N_{\text{samples}}} \|\mathbf{y}_P - \hat{\mathbf{y}}_P\|^2$$

- the resulting training algorithm is called back-propagation (BP) and is essentially an iterative weight update on the basis of a steepest descent criterion.

$$\Delta w_{ji}(t) = -\eta \frac{\partial E}{\partial w_{ji}} + \mu \Delta w_{ji}(t-1)$$

# More on backpropagation

Weight update rule

$$\mathbf{w}(m + 1) = \mathbf{w}(m) + \Delta\mathbf{w}(m),$$

Output neurons

$$\Delta w_{kj} = \eta\delta_k y_j = \eta(t_k - z_k)f'(net_k)y_j.$$

Hidden neurons

$$\Delta w_{ji} = \eta x_i \delta_j = \eta \underbrace{\left[ \sum_{k=1}^{c} w_{kj}\delta_k \right] f'(net_j)}_{\delta_j} x_i.$$

## Second order methods

- Backpropagation is rather simple and easy to implement but can suffer from sever convergence problems.

- One way of coping with this is momentum.

- Another is to use second-order methods.

- Given an initial estimate of the weights, the error is Taylor-expanded up to the second order:

$$E^{\text{batch}}(\mathbf{w}) \cong E^{\text{batch}}(\mathbf{w}_0) + \left(\nabla_{\mathbf{w}_0} E^{\text{batch}}\right)^{\mathbf{T}} \Delta\mathbf{w} + \frac{1}{2}(\Delta\mathbf{w})^{\mathbf{T}} \mathbf{H} \Delta\mathbf{w}$$

- where:

$$\nabla_{\mathbf{w}_0} E^{\text{batch}} = \left.\frac{\partial E^{\text{batch}}}{\partial\mathbf{w}}\right|_{\mathbf{w}=\mathbf{w}_0} \qquad \mathbf{H} = \left.\frac{\partial^2 E^{\text{batch}}}{\partial\mathbf{w}^2}\right|_{\mathbf{w}=\mathbf{w}_0}$$

- Backpropagation corresponds to truncating the expansion to the first order.

# Second order methods - 2

- Based on the second order expansion, solution is sought by differentiating the previous equation wrt to the weights and setting the derivative = 0 (Newton's method)

$$\Delta \mathbf{w}^* = -\mathbf{H}_{\mathbf{w}_0}^{-1} \nabla_{\mathbf{w}_0} E^{\text{batch}}$$

- Problems:
  - Hessian matrix H should be calculated and stored (computationally and memory intensive).
  - Hessian matrix should be nonsingular (not guaranteed, quite often H is not full rank).

- Solutions:
  - Quasi-second order methods:
    - Gauss-Newton
    - Levenberg-Marquardt

## Levenberg-Marquardt method

- In Gauss-Newton, the Hessian matrix is approximated by:

$$H \approx \left( \nabla_{w_0} E^{batch} \right) \left( \nabla_{w_0} E^{batch} \right)^T$$

- However, this approximation is accurate only near the minimum and can suffer from initial guesses far from the optimal solution.

- LM method introduces an additional term to stabilize the estimate:

$$\Delta w_{LM}^* = -H_{w_0}^{-1} \nabla_{w_0} E^{batch} \approx - \left[ \left( \nabla_{w_0} E^{batch} \right) \left( \nabla_{w_0} E^{batch} \right)^T + \mu I \right]^{-1} \left( \nabla_{w_0} E^{batch} \right)$$

- The parameter μ controls the property of the algorithm:
  - When μ is large the steepest descent term is dominant
  - When μ is small the Gauss-Newton term is dominant
- Fast and efficient

# Specifically for classification

- When ANNs are used for classification:
  - Error criterion is cross-entropy instead of RMSE

$$E_{CE} = \sum_{i=1}^{N} \sum_{j=1}^{M} y_{ij} \ln \frac{y_{ij}}{\hat{y}_{ij}}$$

  - Output function is "softmax"

$$\hat{y}_j = \frac{e^{t_j}}{\sum_{k=1}^{M} e^{t_k}}$$

# Generalization *vs* specialization

- Optimal number of hidden neurons
  - Too many hidden neurons: you get an over fit, training set is memorized, thus making the network useless on new data sets
  - Not enough hidden neurons:
    network is unable to learn problem concept

    *~ conceptually: the network's language isn't able to express the problem solution*

# Generalization vs. specialization 2

- Overtraining:
  - Too much examples, the ANN memorizes the examples instead of the general idea

- Generalization vs. specialization trade-off:

  **# hidden nodes & training samples**

# RADIAL BASIS FUNCTIONS - NN

# The RBF-NN

- Differently from MLP, RBF-NN performs classification and regression based on similarity with examples from the training set

- Its basic unit is the Gaussian (RBF) neuron:

$$y = e^{-\frac{\|\mathbf{x}_i - \mathbf{c}_k\|^2}{2\sigma^2}} = e^{-\beta\|\mathbf{x}_i - \mathbf{c}_k\|^2}$$

$c_k$ is the center of the $k^{th}$ RBF

# The RBF-NN for classification

$$y_j = \sum_{k=1}^{n_{RBF}} w_{jk} e^{-\beta \|\mathbf{x}_i - \mu_k\|^2}$$

**RBF Neurons**

**Input Vector**

**Weighted Sums**

$w_{11}$  Cat. 1 Weights

$\Sigma$  →  Category 1 Score  $y_1$

$\mu_1$

$\mu_2$

$\Sigma$  →  Category c Score  $y_c$

Cat. c Weights

$w_{ck}$

$\mu_k$

$\mathbf{x}_i$

$\mu$ is the prototype to compare against

There is the need to optimize the center and width of RBFs and the weights $w_{jk}$

# Training the RBF-NN

- There are many different training algorithms for RBF-NN
- Apart from backpropagation, the most common are Orthogonal least squares and the following:

  1. Select the RBF centers by k-means clustering (possibly, applying clustering separately by category)
  2. Calculate the RBF width as the mean cluster distance to the centroid:

$$\sigma_k = \frac{1}{m} \sum_{i=1}^{m} \|\mathbf{x}_i - \mu_k\| \quad \text{and} \quad \beta_k = \frac{1}{2\sigma_k^2}$$

  1. Calculate the weight by LS regression:

$$\mathbf{W} = \mathbf{H}^{\dagger}\mathbf{Y} \quad \text{with } h_{ik} = e^{-\beta_k \|\mathbf{x}_i - \mu_k\|^2}$$

# The RBF-NN for regression

$$y = \overset{\circ}{\text{a}}\,_{k=1}^{n_{RBF}} w_k e^{-\beta\|\mathbf{x}_i - \mu_k\|^2}$$

**RBF Neurons**

**Input Vector**

$w_1$

$\mu_1$

$\mu_2$

$w_k$

$\mu_k$

$\mathbf{x}_i$

Output Weights

**Output Node**

$\Sigma$

Function Output

$y$

$\mu_k$ = Prototype for neuron k

There is the need to optimize the center and width of RBFs and the weights $w_{jk}$

# Training the RBF-NN for regression

- The training algorithms for regression are the same as for classification, with two main differences:
  1. The width of the RBF is normally selected as equal for all nodes and based on cross-validation instead as on clustering.
  2. The output of the function is often scaled (normalized):

$$y = \frac{\sum_{k=1}^{n_{RBF}} w_k e^{-\beta \|\mathbf{x}_i - \mu_k\|^2}}{\sum_{k=1}^{n_{RBF}} e^{-\beta \|\mathbf{x}_i - \mu_k\|^2}}$$

# Self organizing maps (Kohonen architectures)

# Relative distance measure

The objects $X_s$ are adjusted according to an absolute distance measure $d(X_s, W_j)$ to positions of the prespecified points $W_j$ distributed in a topologically predefined scheme.

In the case of Kohonen neural networks the objects $X_s$ and the arbitrarily distributed points $W_j$ are, first, assigned to each other, and next, the points $W_j$ with associated closest objects $\{X_s\}$ are pooled together to predefined positions in a 2-D plane.



**W**

$X_s$

W'$_j$

**The final result does not depend much on the distances between objects, but rather on the distances between the objects $X_s$ and the pivot points $W_j$.**

# Kohonen self organizing maps (SOMs)

- The implicit functional relation that we want to approximate is a nonlinear mapping from an $N_i$-dimensional input space to a low-dimensional (usually 2D) discrete coordinate space (the map).

- Since there is no desired response to be obtained, training occurs by self-organization, i.e. a Kohonen network adapts itself so that similar input objects are associated with topological close neurons.

- In a self-organizing map, the target space used in Kohonen mapping is a two-dimensional array of neurons (the Kohonen layer or top-map), fully connected to the input layer, onto which the samples are mapped.

- Introducing the preservation of topology results in specifying for each node in the Kohonen layer a defined number of neurons as nearest neighbors, second-nearest neighbors and so on.

# Kohonen self organizing maps 2

The most important feature of the Kohonen neural network is the topological order in which the neurons are combined together into the network.

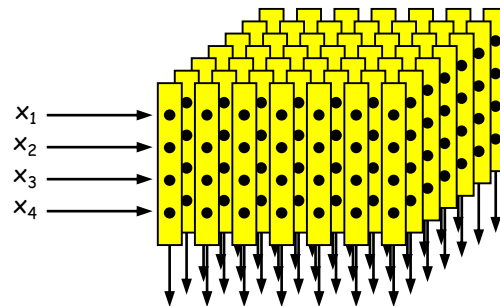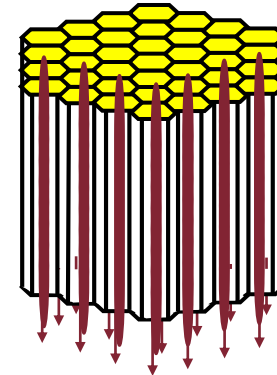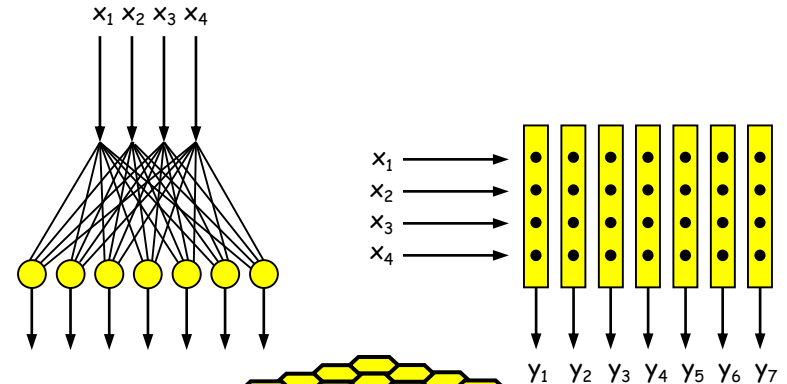The lay-out of neurons in the Kohonen network can be linear (all neurons in one row ot line) or planar (in a rectangular or hexagonal lay-out).

$x_1$ $x_2$ $x_3$ $x_4$

$x_1$
$x_2$
$x_3$
$x_4$

$y_1$ $y_2$ $y_3$ $y_4$ $y_5$ $y_6$ $y_7$

$x_1$ $x_2$ $x_3$ $x_4$

$y_{11}$ $y_{12}$ $y_{13}$ $y_{14}$ $y_{15}$ $y_{16}$ $y_{17}$
$y_{21}$ $y_{22}$ $y_{23}$ $y_{24}$ $y_{25}$ $y_{26}$ $y_{17}$
$y_{31}$ $y_{32}$ $y_{33}$ $y_{34}$ $y_{35}$ $y_{36}$ $y_{37}$
$y_{41}$ $y_{42}$ $y_{43}$ $y_{44}$ $y_{45}$ $y_{46}$ $y_{47}$
$y_{51}$ $y_{52}$ $y_{53}$ $y_{54}$ $y_{55}$ $y_{56}$ $y_{57}$
$y_{16}$ $y_{62}$ $y_{63}$ $y_{64}$ $y_{65}$ $y_{66}$ $y_{67}$
$y_{17}$ $y_{72}$ $y_{73}$ $y_{74}$ $y_{75}$ $y_{76}$ $y_{77}$

$x_1$
$x_2$
$x_3$
$x_4$

$y_{11}$ $y_{12}$ $y_{13}$ $y_{14}$ $y_{15}$ $y_{16}$ $y_{17}$
$y_{21}$ $y_{22}$ $y_{23}$ $y_{24}$ $y_{25}$ $y_{26}$ $y_{17}$
$y_{31}$ $y_{32}$ $y_{33}$ $y_{34}$ $y_{35}$ $y_{36}$ $y_{37}$
$y_{41}$ $y_{42}$ $y_{43}$ $y_{44}$ $y_{45}$ $y_{46}$ $y_{47}$
$y_{51}$ $y_{52}$ $y_{53}$ $y_{54}$ $y_{55}$ $y_{56}$ $y_{57}$
$y_{16}$ $y_{62}$ $y_{63}$ $y_{64}$ $y_{65}$ $y_{66}$ $y_{67}$
$y_{17}$ $y_{72}$ $y_{73}$ $y_{74}$ $y_{75}$ $y_{76}$ $y_{77}$

$y_{74}$
$y_{63}$ $y_{75}$
$y_{52}$ $y_{64}$ $y_{76}$
$y_{41}$ $y_{53}$ $y_{65}$ $y_{77}$
$y_{42}$ $y_{54}$ $y_{66}$
$y_{31}$ $y_{43}$ $y_{55}$ $y_{67}$
$y_{32}$ $y_{44}$ $y_{56}$
$y_{21}$ $y_{33}$ $y_{45}$ $y_{57}$
$y_{22}$ $y_{34}$ $y_{46}$
$y_{11}$ $y_{23}$ $y_{35}$ $y_{47}$
$y_{12}$ $y_{24}$ $y_{36}$
$y_{13}$ $y_{25}$
$y_{14}$

# Defining the neighborhood

The *neighborhood* of a neuron is usually considered to be square or hexagonal which means that each neuron has 8 or 6 nearest neighbors respectively.

# Cyclic and toroid conditions

The cyclic and the toroid boundary conditions. The edge one one side is linked to the edge on the opposite side.

Cyclic conditions in a line of neurons: 1st neighbours of $W_1$ are $W_2$ and $W_n$

0  1  2  3 …        … 3  2  1  0

Toroid conditions in a rectangular plane of neurons: 1st neghbours of neurons on edges $a$ and $c$ are neurons on edges $b$ and $d$ respectively.

# Cyclic and toroid conditions 2

The consideration of the toroidal conditions in Kohonen neural networks can sometimes lead to much clearer results.

# Learning procedure in Kohonen networks

Learning in the Kohonen neural network is iterative, i.e., a set of objects $\{X_s\}$ is sent through the network several times.

After the pass of one object through the network the weights, which at the begining of learning are randomised, are changed. The learning procedure of a single pass consists of three steps:

1. **selection** of one neuron in the network according to a prespecified criterion   (largest response, most similar input $X_s$ to the neuron $W_j$, or similar),

2. **correction of the weights of the selected neuron**, and

3. **correction of the weights of the neighbouring** neurons up to a specified range  around the selected neuron.

One pass of the entire data set through the network is called one epoch.

Training  is stopped after a certain number of epochs. The main result of the Kohonen learning is the "top-map" or "self-organised map" (SOM) of objects associated with the excited neurons.

# Kohonen networks in practice

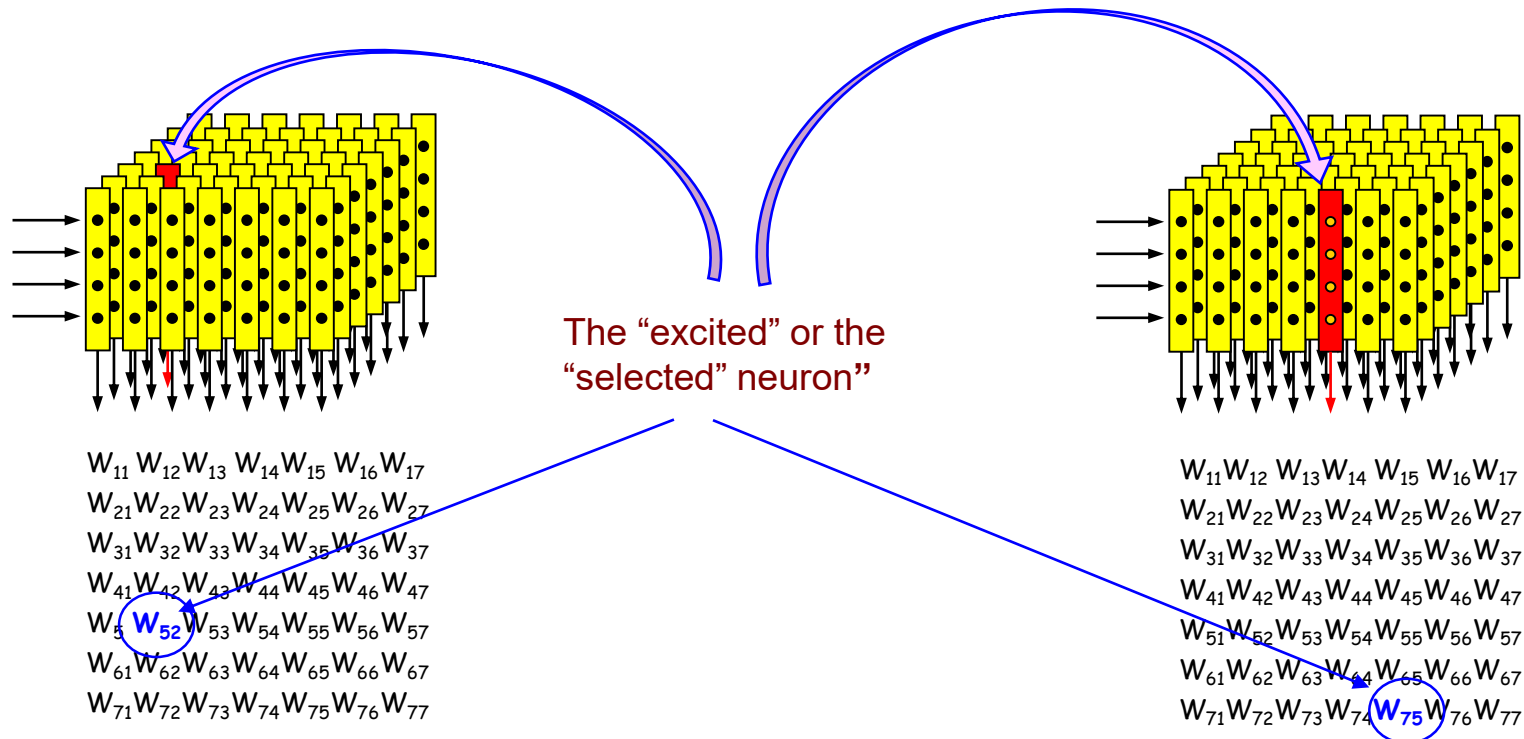**Step one:** selection of the "excited", "central", or "responding" neuron. In large majority of cases the selection of the neuron is made according to the smallest distance criterion:

$$d(X_s, W_j) = \sum_{i=1}^{m} (x_{si} - w_{ji})^2 \quad \text{for all } j = 1, \ldots, N \times N$$

The "excited" or the "selected" neuron"

$W_{11}\ W_{12}\ W_{13}\ W_{14}\ W_{15}\ W_{16}\ W_{17}$
$W_{21}\ W_{22}\ W_{23}\ W_{24}\ W_{25}\ W_{26}\ W_{27}$
$W_{31}\ W_{32}\ W_{33}\ W_{34}\ W_{35}\ W_{36}\ W_{37}$
$W_{41}\ W_{42}\ W_{43}\ W_{44}\ W_{45}\ W_{46}\ W_{47}$
$W_{51}\ \mathbf{W_{52}}\ W_{53}\ W_{54}\ W_{55}\ W_{56}\ W_{57}$
$W_{61}\ W_{62}\ W_{63}\ W_{64}\ W_{65}\ W_{66}\ W_{67}$
$W_{71}\ W_{72}\ W_{73}\ W_{74}\ W_{75}\ W_{76}\ W_{77}$

$W_{11}\ W_{12}\ W_{13}\ W_{14}\ W_{15}\ W_{16}\ W_{17}$
$W_{21}\ W_{22}\ W_{23}\ W_{24}\ W_{25}\ W_{26}\ W_{27}$
$W_{31}\ W_{32}\ W_{33}\ W_{34}\ W_{35}\ W_{36}\ W_{37}$
$W_{41}\ W_{42}\ W_{43}\ W_{44}\ W_{45}\ W_{46}\ W_{47}$
$W_{51}\ W_{52}\ W_{53}\ W_{54}\ W_{55}\ W_{56}\ W_{57}$
$W_{61}\ W_{62}\ W_{63}\ W_{64}\ W_{65}\ W_{66}\ W_{67}$
$W_{71}\ W_{72}\ W_{73}\ W_{74}\ \mathbf{W_{75}}\ W_{76}\ W_{77}$

# Kohonen networks in practice 2

Step two: correction of weights of the selected neuron $W_e$

Because after each pass only one neuron is selected the leaning procedure is called the "winner-takes-all" strategy.

$$\otimes W_e = \eta \, (X_s - W_e^{old})$$

$$d \, (W_e^{new}, X_s)$$

$$W_e$$

$$W_e^{new}$$

$$X_s$$

$$\eta = (\eta^{start} - \eta^{final})(1 - n_{epoch}/n_{tot}) + \eta^{final}$$
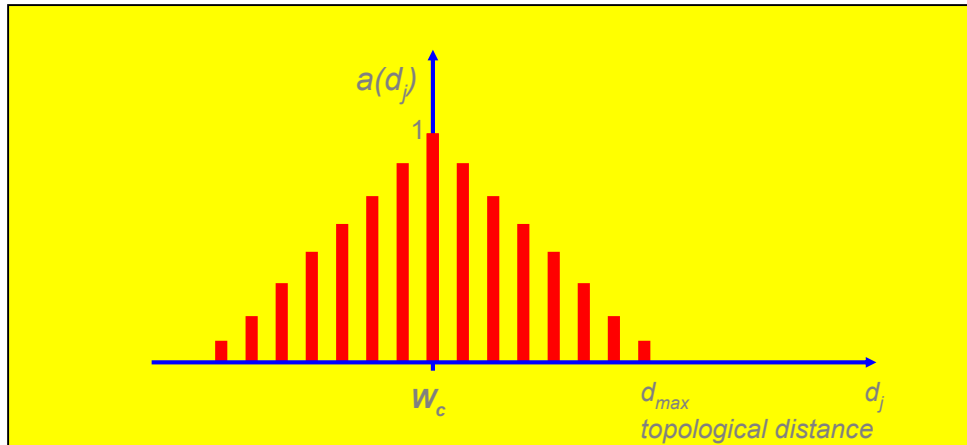
The parameter $\eta$ is called the learning rate and is in most applications "time" dependent, i.e. dependent on the number of currently performed learning epochs.

The corrections $\otimes W_e$ are driving the weights $w_{ei}$ of the excited neuron $W_e = (w_{e1}, w_{e2}, ... w_{ei}...w_{em})$ closer to the variables $x_{si}$ of object $X_s = (x_{s1}, x_{s2}, ... x_{si}, ... x_{sm})$ that has excited it.

# Kohonen networks in practice 3

Step three: correction of the weights in neurons $W_j$ surrounding the selected neuron $W_e$

$$\Delta w_{ji} = \eta a(d_j)\left(x_{si} - w_{ji}^{old}\right)$$

$$a(d_j) = 1 - \frac{d_j}{(d_{max}+1)}$$



For the selected neuron $a(d_j) = 1$ because $d_j = 0$.
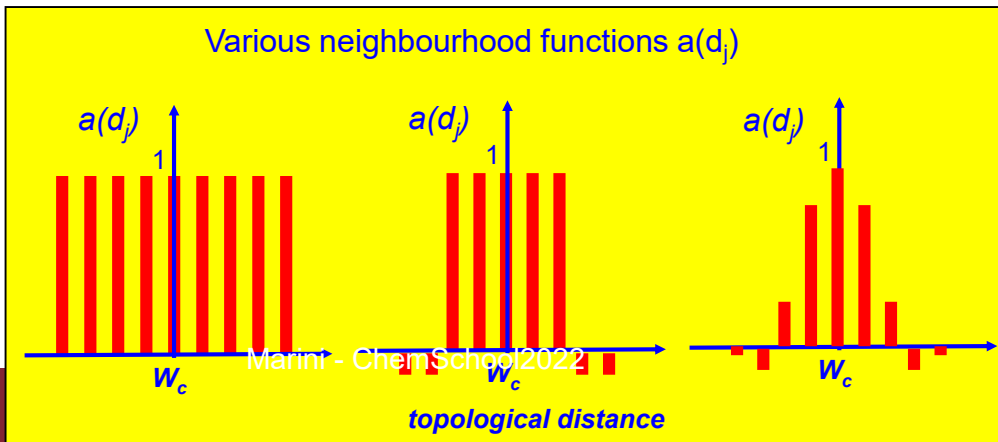
For the neurons separated $d_{max}$ from the selected neuron $W_e$ the value $a(d_{max})$ is the smallest posible corection.

$$\Delta w_{ji} = \eta(1 - \frac{d_j}{d_{max}+1})(x_{si} - w_{ji}^{old})$$

$$\Delta w_{ji} = \eta(1 - \frac{d_j}{N_{net}(1 - \frac{n_{epoch}}{n_{tot}})+1})(x_{si} - w_{ji}^{old})$$
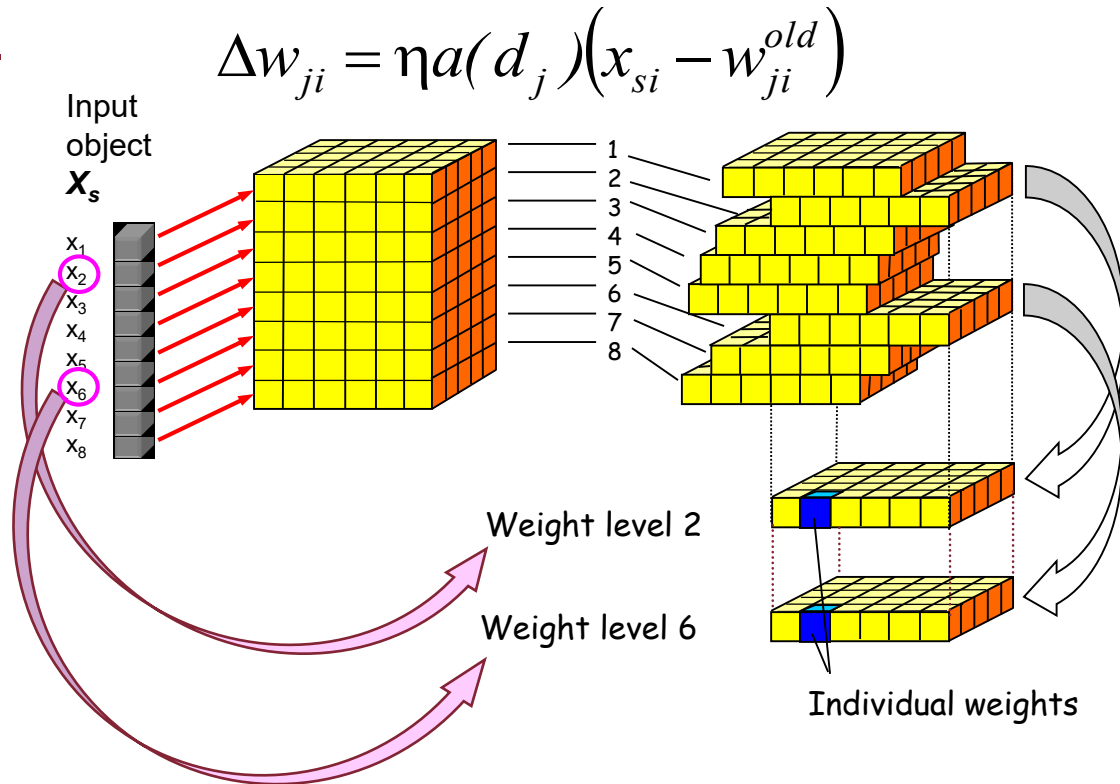
At the beginning of learning:
$n_{epoch} = 1$  ➔  $d_{max} = N_{net}$; correction covers the entire network

At the end of the learning:
$n_{epoch} = n_{tot}$  ➔  $d_{max} = 0$; correction is limited only to the selected neuron $W_e$.



Various neighbourhood functions $a(d_j)$

topological distance

# Weight maps

Correction of the weights in neurons $W_j$ is applied to all levels of weights. Due to the fact that weights $w_{ji}$ in the neurons of the Kohonen network are alligned into levels according to the order of input variables $x_i$, each level of weights is influenced only by one variable and thus forms a map of weight values or weight map.

$$\Delta w_{ji} = \eta a(d_j)\left(x_{si} - w_{ji}^{old}\right)$$

Input
object
$X_s$

$x_1$
$x_2$
$x_3$
$x_4$
$x_5$
$x_6$
$x_7$
$x_8$

1
2
3
4
5
6
7
8

Weight level 2

Weight level 6

Individual weights

# Weight maps 2

| | | | | | | |
|---|---|---|---|---|---|---|
| A | A | A | A,C | | | B |
| A | A | | A,C | C | | B |
| A | | C | C | C | | B |
| | C | C | C | | | B |
| C | C | C | | | B | B |
| C | C | C,B | B | B | B | B |
| C | C | | B | B | B | B |

```
.76 .85 .90 .93 .95 .97 .99

.63 .75 .84 .80 .85 .89 .95

.59 .66 .68 .65 .77 .75 .92

.32 .33 .52 .58 .60 .69 .85

.31  43 .46 .42 .53 .61 .81

.16 .25 .30 .33 .46 .55 .75

.03 .09 .08 .22 .38 .50 .67
```

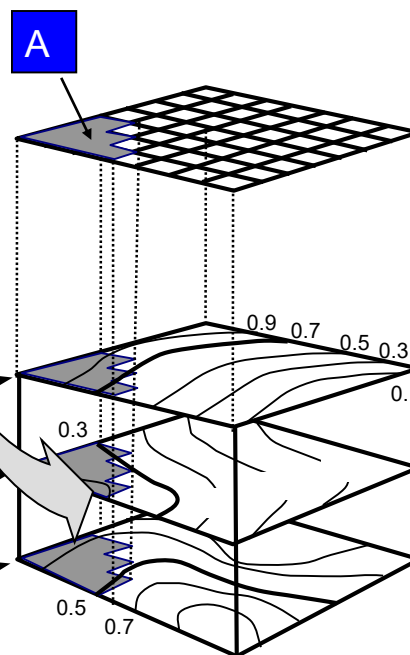Weight map of the variable $x_2$ - *normalised binder concentration*

A

Input object $\boldsymbol{X_s}$ (recipe for the paint-coat).  Only three variables are shown.

Solvent     $x_1$

Binder      $x_2$

Pigment    $x$

0.9  0.7  0.5  0.3  0.1
0.3
0.5  0.7

Top map showing the cells containing the information avout the overal quality of the paint-coat

Labels **A**, **B**, and **C** refer to the excellent, passable, and below standard quality, respectively.

The most important feature of the Kohonen network is the fact that all weights in all neurons regardless whether they were excited during the trainng or not bear valuable information.
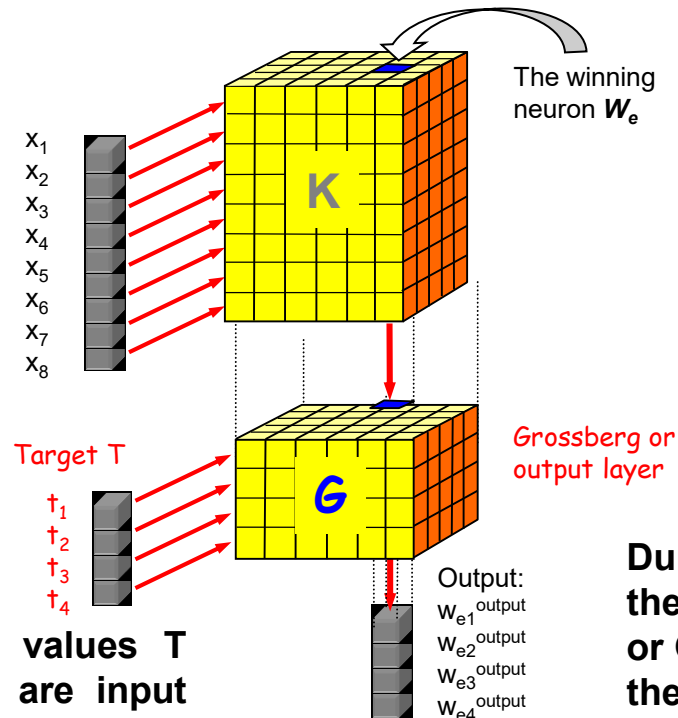
# Counterpropagation networks

By the addition of an identical layer of weights to which the targets are input the Kohonen network is transformed into the counterpropagation network.

The additional layer (the output or Grossberg layer) has the same number and the same layout of neurons as the first (Kohonen) layer, however, in each layer the neurons have different number of weights.

$$\Delta w_{ji}^{K} = \eta a(d_j)\left(x_{si} - w_{ji}^{K,old}\right)$$

$x_1$
$x_2$
$x_3$
$x_4$
$x_5$
$x_6$
$x_7$
$x_8$

**K**

The winning neuron $W_e$

$$\Delta w_{ji}^{G} = \eta a(d_j)\left(t_{si} - w_{ji}^{G,old}\right)$$

Target T

$t_1$
$t_2$
$t_3$
$t_4$

**G**

Grossberg or output layer

Output:
$w_{e1}^{output}$
$w_{e2}^{output}$
$w_{e3}^{output}$
$w_{e4}^{output}$

**During the training the target values T associated with each object X are input into the output layer in exactly the same manner as the objects X are input into the Kohonen layer.**
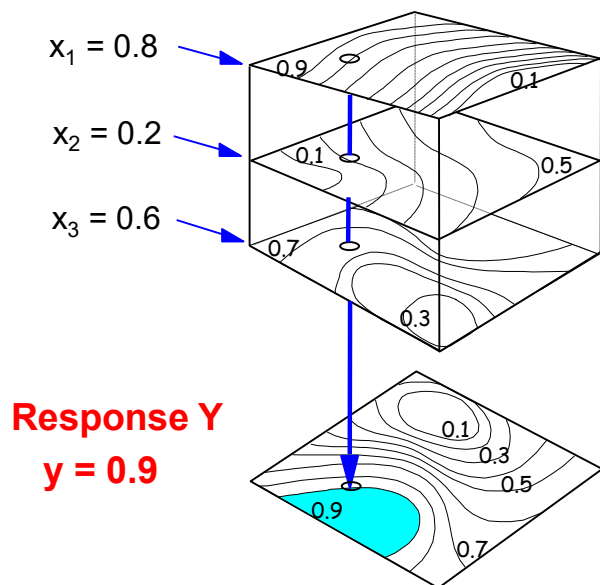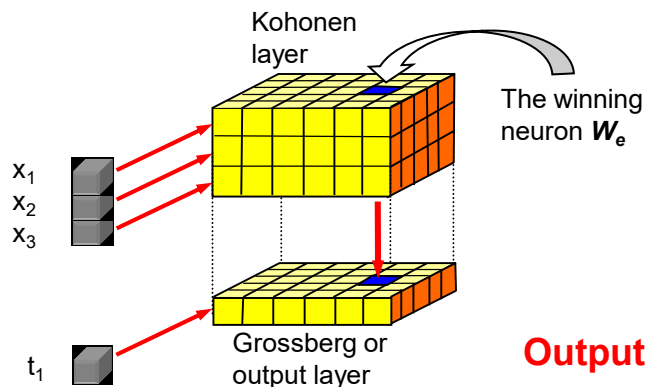
**During the retrieval (prediction) the weights $w_{ei}^{output}$ in the output or Grossberg layer selected by the winning neuron of the input object X (in the Kohonen layer) are used as the predictions.**

# Counterpropagation networks 2

Each counterpropagation network can be used as a direct and as an inverse "model".



Kohonen layer

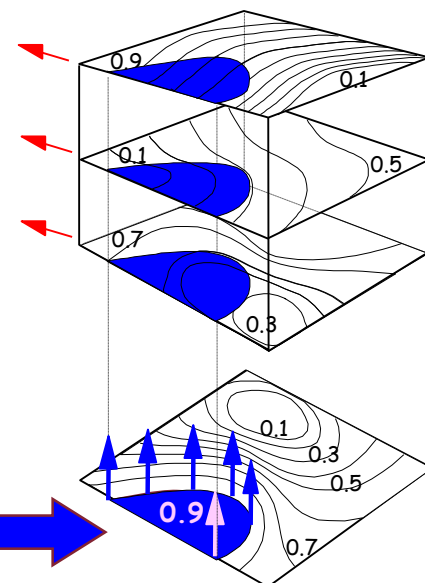The winning neuron $W_e$

$x_1$
$x_2$
$x_3$

$t_1$

Grossberg or output layer

**Output X**

$x_1 = (0.45 - 0.95)$

$x_2 = (0.00 - 0.30)$

$x_3 = (0.30 - 0.70)$

$x_1 = 0.8$

$x_2 = 0.2$

$x_3 = 0.6$

**Response Y**
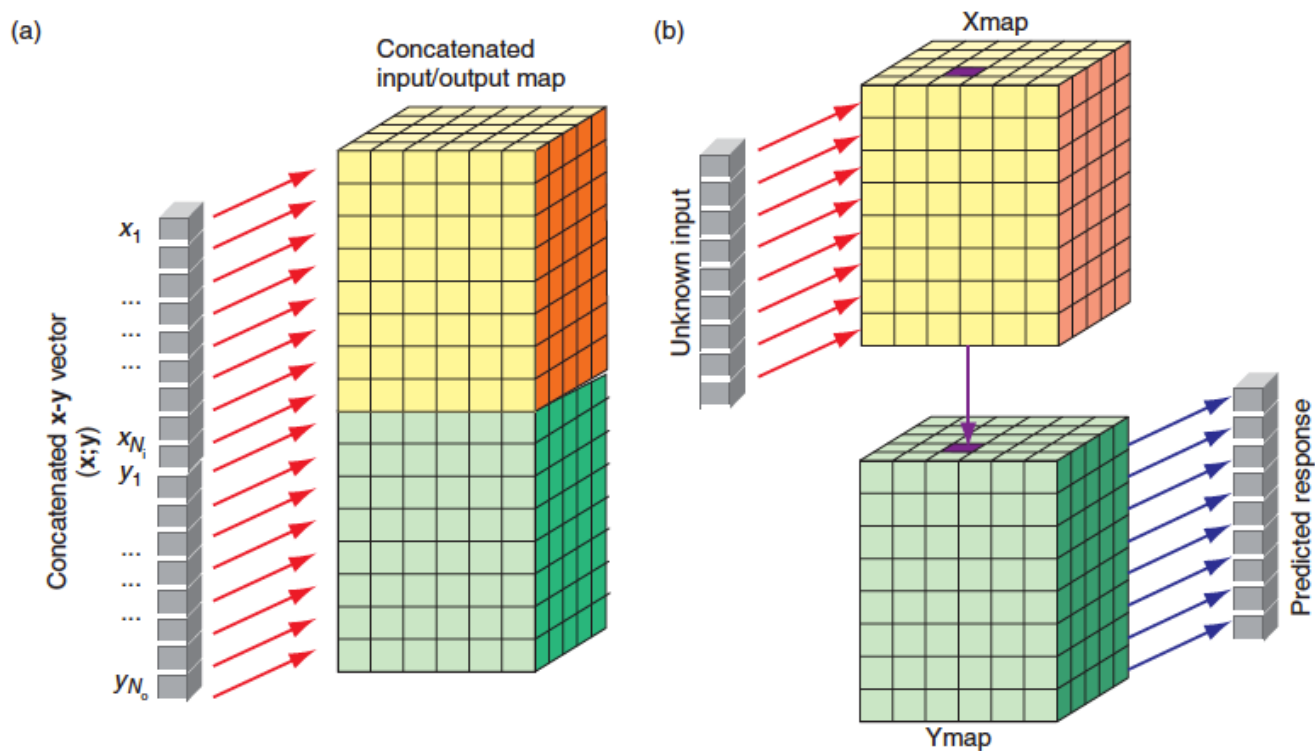**y = 0.9**

Input Y

y = 0.9

# Supervised Kohonen architectures

- Counterpropagation networks are semi-supervised architectures, as the value of the Y vector doesn't drive the selection of the winning neuron and, as a consequence, the direction of the training.

- There are other proposed architectures that are truly supervised and that can be used to build classification models:

  - Supervised Kohonen networks

  - XY-fused network

  - Bidirectional Kohonen networks

# Supervised Kohonen networks

- X and Y variables are concatenated to train the network as in the standard Kohonen architecture.

- After training the two blocks are separated and prediction occur as in counterpropagation

# XY-fused networks

- The winning neuron is decided by considering a similarity function which is a weighted sum of the similarity in the X space and in the Y space.

- The parameter α starts with a high value and then decreases, so that at the end only similarity in the Y space contributes, while at the beginning the X space is dominant.

$$S_{\text{fused}}(j) = \alpha(t)S(\mathbf{x}, \mathbf{w}_j) + (1 - \alpha(t))S(\mathbf{y}, \mathbf{u}_j)$$

# Bidirectional Kohonen networks

- The concept of BDK networks is similar to that of XY-fused, combining similarity in both spaces to update the weights.

- In BDK, however, X and Y weights are updated in an alternate fashion:

- The parameter α starts with a high value and then decreases, so in the beginning similarity in the Y space governs selection of the winning neuron in the X space while similarity in the X spaces de

$$S_{\text{Winner}\,X}(j) = (1 - \alpha(t))\,S(\mathbf{x}, \mathbf{w}_j) + \alpha(t)\,S(\mathbf{y}, \mathbf{u}_j)$$

$$S_{\text{Winner}\,Y}(j) = \alpha(t)\,S(\mathbf{x}, \mathbf{w}_j) + (1 - \alpha(t))\,S(\mathbf{y}, \mathbf{u}_j)$$

# The fortunes and misfortunes of NNs

- Despite their increasing popularity up to the beginning of the 2000s, interest in neural networks seemed to vanish:
    - Curse of dimensionality
    - Inefficient learning algorithms
    - Lack of interpretability of the models
    - Performances heavily dependent on the choice of data representation (or features) on which they are applied.

- Much of the actual effort in deploying machine learning algorithms goes into the design of preprocessing pipelines and data transformations that result in a representation of the data that can support effective machine learning

# Representation learning

- Learning representations of the data that make it easier to extract useful information when building classifiers or other predictors
  - Captures the posterior distribution of the underlying explanatory factors for the observed input.
  - Is useful as input to a supervised predictor.
- DEEP LEARNING:
  - composition of multiple non-linear transformations, with the goal of yielding more abstract – and ultimately more useful – representations
- Fundamental questions:
  - What makes one representation better than another?
  - Given an example, how should we compute its representation, i.e. perform feature extraction?
  - What are appropriate objectives for learning good representations?

# Representation learning: "chemometric" concepts

- Representations → convenient to express many general (not task-specific) priors that are likely to be useful for a learning machine to solve AI-tasks.

- The revival experienced by NNs in the recent years has much to do with such priors
  - Their absence was one of the main reasons for the vanishing interest towards NN in the 2000s
  - They share much in common with essential ideas which make chemometric representations useful and versatile
  - Possibility of mutual benefit between the disciplines

- Some of these will be briefly discussed in the following

# The general priors of representation learning

- **Smoothness**:
    - $x_1 \approx x_2 \Rightarrow f(x_1) \approx f(x_2)$
    - Insufficient to circumvent the curse of dimensionality
- **Multiple explanatory factors**:
    - Data distribution generated by different underlying factors
    - What one learns about one factor generalizes in many configurations of the other factors
    - The objective is to recover or at least disentangle these underlying factors of variation
- **Hierarchy**:
    - The features that are useful for describing the world around us can be defined in terms of other features, in a hierarchy
    - More abstract concepts higher in the hierarchy are defined in terms of less abstract ones → **Deep learning**

# The general priors of representation learning - 2

- **Semi-supervised learning**:
    - A subset of the factors explaining X's distribution explain much of Y, given X.
    - Representations that are useful for P(X) tend to be useful when learning P(Y|X)
    - Sharing of statistical strength between the unsupervised and supervised learning tasks

- **Manifolds**:
    - Probability mass concentrates near regions that have a much smaller dimensionality than the original space where the data lives.

- **Sparsity**:
    - For any given x, only a small fraction of the possible factors are relevant.
    - Features that are often zero or insensitive to small variations of x.
    - Priors on latent variables (peaked at 0), or by a nonlinearity whose value is often flat at 0 (e.g., ReLU)
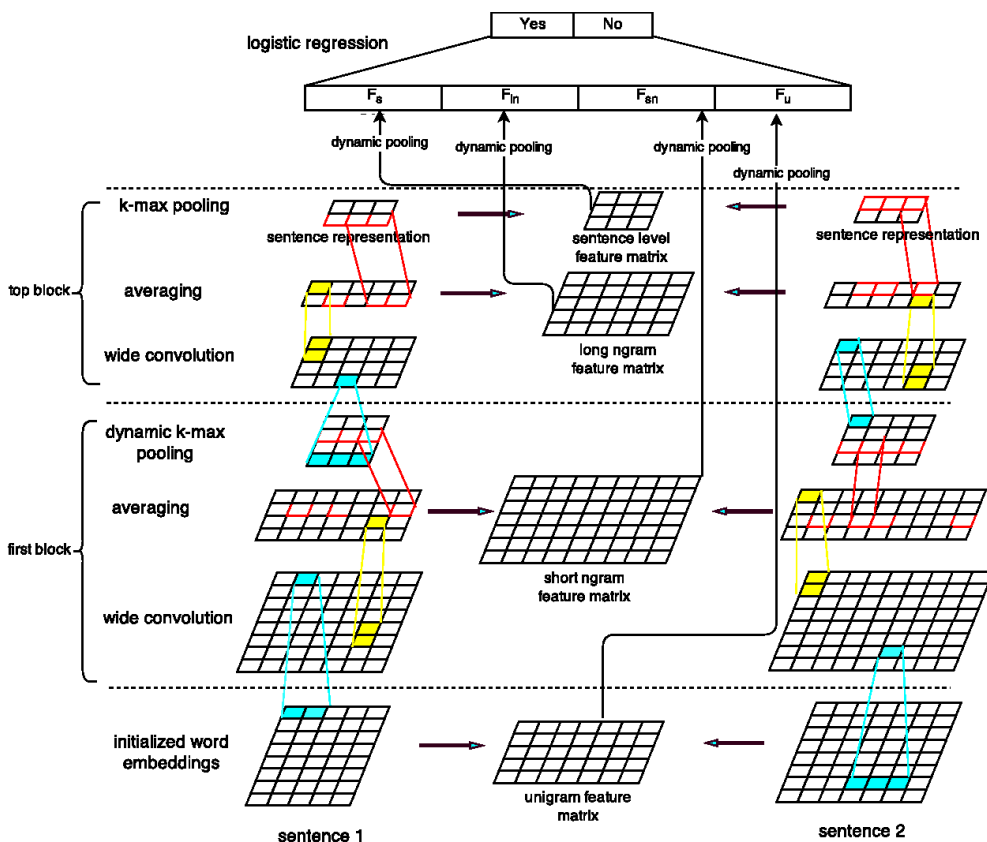
# Building deep representations



- Learn a hierarchy of features one level at a time, using unsupervised feature learning to learn a new transformation at each level to be composed with the previously learned transformations;
- The set of layers could be combined to initialize a deep supervised predictor, such as a neural network classifier, or a deep Boltzmann Machine
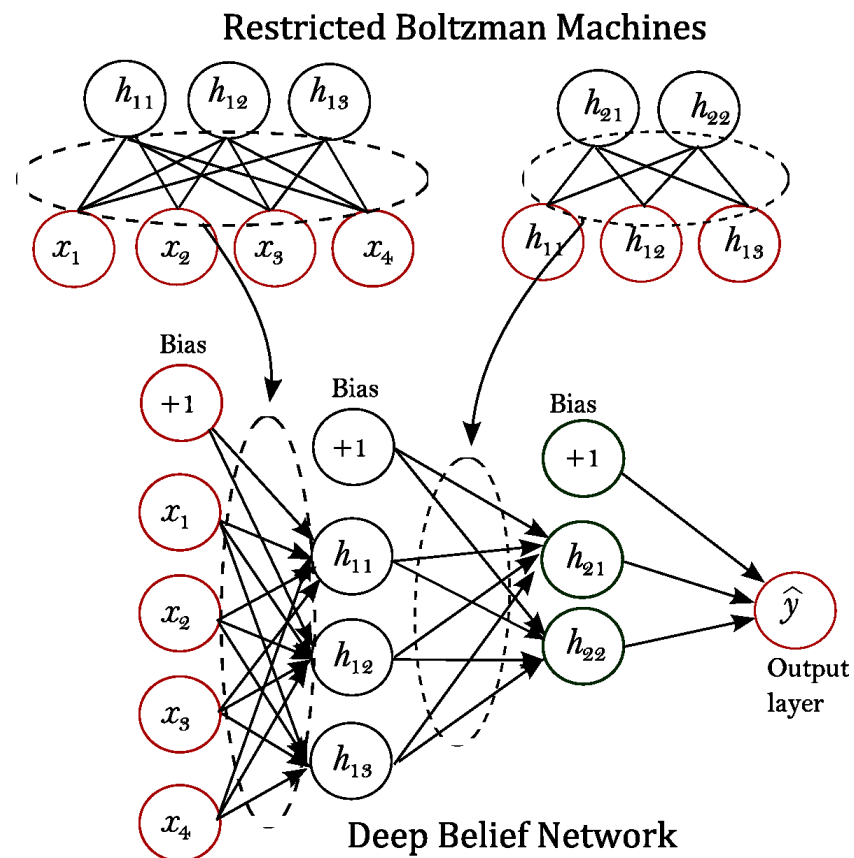
# Two learning paradigms

- One rooted in probabilistic graphical models and one rooted in neural networks.

- Probabilistic modeling:
  - Attempt to recover a parsimonious set of latent random variables that describe a distribution over the observed data.
  - Feature values are conceived as the result of an inference process to determine the probability distribution of the latent variables given the data, i.e. $p(h|x)$, the posterior probability.

- Main differences:
  - The layered architecture of a deep learning model is to be interpreted as describing a probabilistic graphical model or as describing a computation graph?
  - Are hidden units considered latent random variables or as computational nodes?
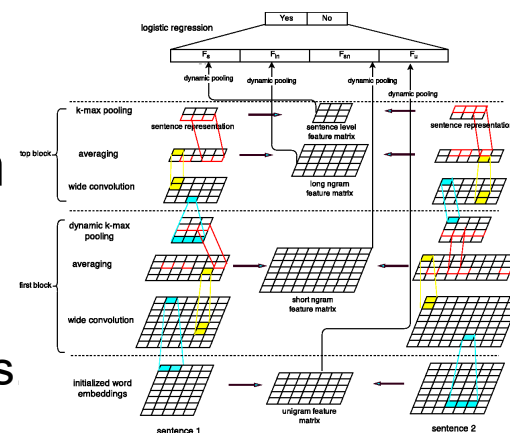
# And two corresponding architectures



- Convolutional NN
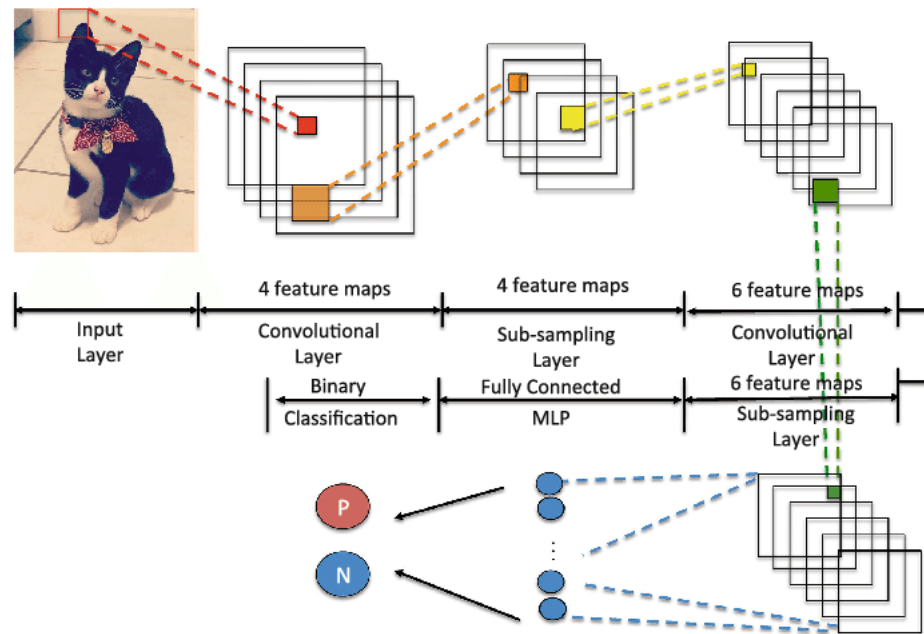
- (Restricted) Boltzmann machines

# Convolutional Neural Networks



- Designed to process data that come in the form of multiple arrays:
  - a colour image composed of three 2D arrays containing pixel intensities in the three colour channels
  - 1D for signals and sequences

- There are four key ideas behind ConvNets that take advantage of the properties of natural signals:
  - local connections
  - shared weights
  - pooling
  - the use of many layers.

# Convolutional Neural Networks

- The architecture is structured as a series of stages.

- The first few stages are composed of two types of layers: convolutional layers and pooling layers.

    – Units in a convolutional layer are organized in feature maps, within which each unit is connected to local patches in the feature maps of the previous layer through a set of weights called a filter bank.

# The convolutional Concept

- Convolution extracts feature from the input image (data)
- Preserves the spatial relationship between pixels



Image

Convolved Feature

| Operation | Filter | Convolved Image |
|---|---|---|
| Identity | $\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$ |  |
| Edge detection | $\begin{bmatrix} 1 & 0 & -1 \\ 0 & 0 & 0 \\ -1 & 0 & 1 \end{bmatrix}$ |  |
| | $\begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$ |  |
| | $\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$ |  |
| Sharpen | $\begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$ |  |
| Box blur (normalized) | $\frac{1}{9}\begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$ |  |
| Gaussian blur (approximation) | $\frac{1}{16}\begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$ |  |

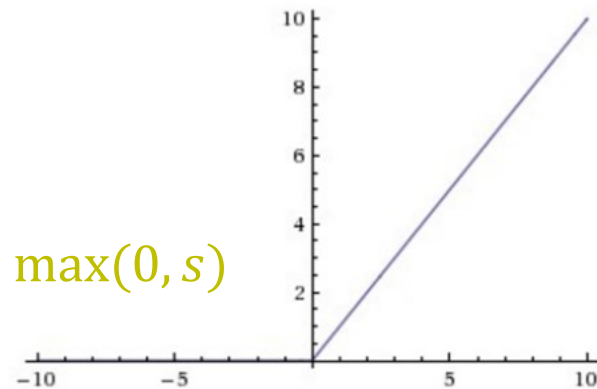# And for multiple-layered inputs

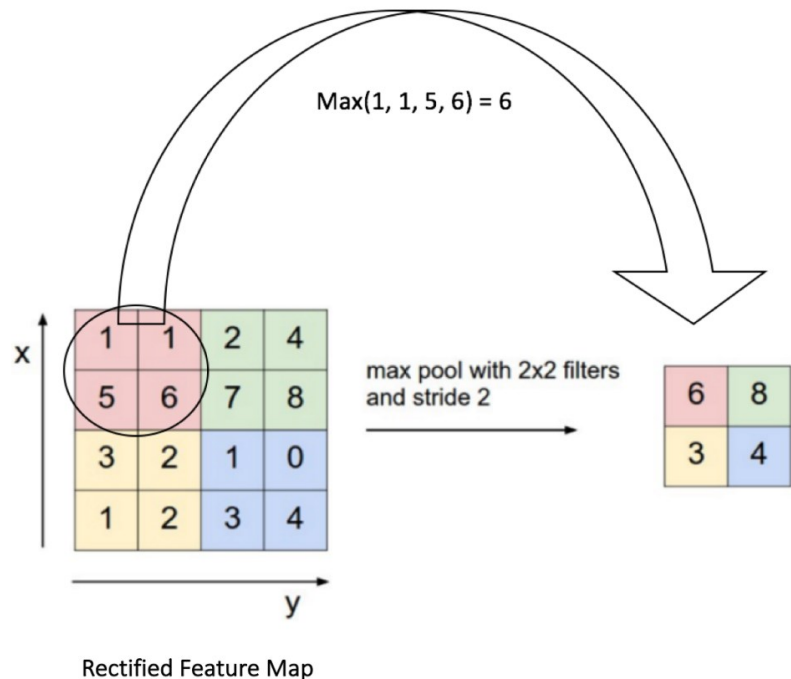- Filtering proceeds in parallel across the depth of the image

# Going nonlinear

- A nonlinear operation is added on top of the convolution.
- Usually it is carried out by means of a Rectified Linear Unit (but one could also use, sigmoid or hyperbolic tangent)
- ReLU is an element wise operation (applied per pixel to the activation maps) and replaces all negative pixel values in the feature map by zero.

$$h = \max(0, s)$$
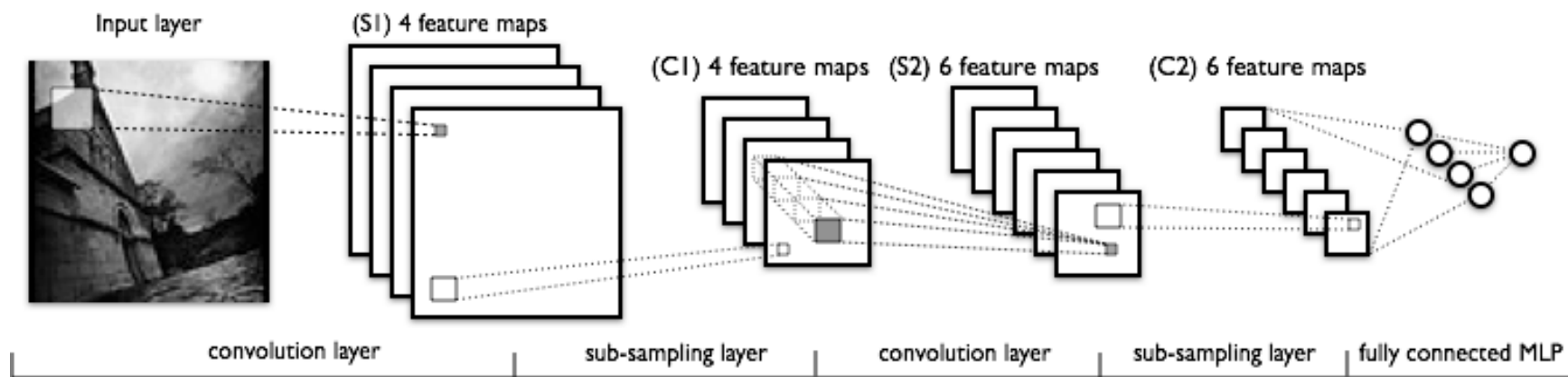
# Compressing (Pooling)

- Spatial Pooling reduces the dimensionality of each feature map but retains the most important information.

- Can be of different types: Max, Average, Sum etc.

- Define a spatial neighborhood (for example, a 2 × 2 window) and:
  – take the largest element from the rectified feature map within that window (Max pooling)
  – Take the average (Average Pooling) or sum of all elements in that window
  – Max Pooling has been shown to work better.

Max(1, 1, 5, 6) = 6

| 1 | 1 | 2 | 4 |
|---|---|---|---|
| 5 | 6 | 7 | 8 |
| 3 | 2 | 1 | 0 |
| 1 | 2 | 3 | 4 |

x

y

max pool with 2x2 filters
and stride 2

| 6 | 8 |
|---|---|
| 3 | 4 |

Rectified Feature Map

# Compressing (Pooling) - 2

- Pooling is applied separately to each feature map
- The function of Pooling is to progressively reduce the spatial size of the input representation:
  - makes the input representations (feature dimension) smaller and more manageable
  - reduces the number of parameters and computations in the network, therefore, controlling overfitting
  - makes the network invariant to small transformations, distortions and translations in the input image (a small distortion in input will not change the output of Pooling – since we take the maximum / average value in a local neighborhood).
  - almost scale invariant representation of our image (the exact term is "equivariant")→detect objects in an image no matter where they are located.

# Wrapping up

# Generative Topographic Mapping (GTM)
## *Eric Latrille – LBE –INRAE (France*

GTM is a dimensionality reduction algorithm well described by Bishop *et al*.

Briefly speaking, the algorithm injects a 2D hypersurface (*manifold*) into an initial *D*dimensional data space. The manifold is fitted to the data distribution by the ExpectationMaximization (*EM*) algorithm which minimizes the log-likelihood of the training data.

Once the fitting is done, each item from the data space is projected to a 2D latent grid of *K* nodes.

GTM is a probabilistic extension of SOM where log-likelihood is utilized as an objective function.

The manifold used to bind a data point **t\*** in the data space and its projection **x\*** in the latent space is described by a set of *M* Radial Basis Function centers (*RBF*; Gaussian functions are generally used).