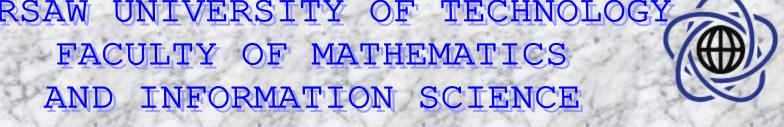
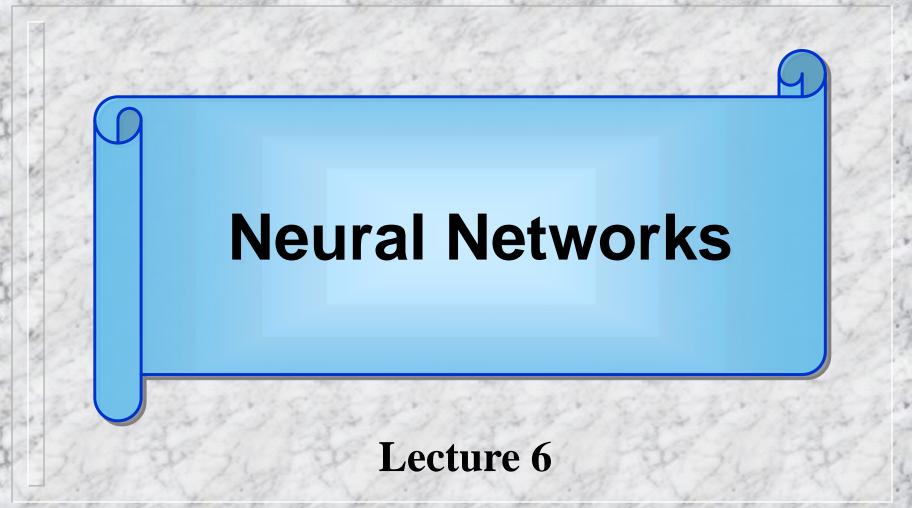


WARSAW UNIVERSITY OF TECHNOLOGY FACULTY OF MATHEMATICS



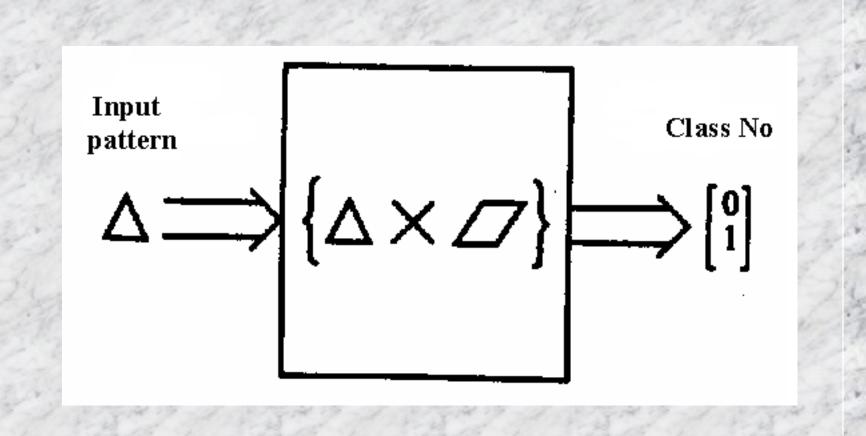




The fundamental objective for pattern recognition is classification and it is the most typical form of neural transformation. A pattern recognition system can be treated as a two stage device: feature extraction and classification. A feature is defined as a measurement taken on the input pattern that is to be classified but the classifier has to map the input features onto a classification state. The classifier must to decide which type of class category they match most closely.

Any given input pattern must belong to one of the classes that are in consideration. So, the mapping from the input to the required class must exists.

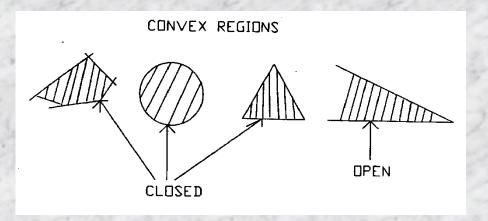
This mapping is a function that transforms the input pattern into the correct output class, and we will consider that our network has learnt to perform correctly, if it can carry out this mapping.



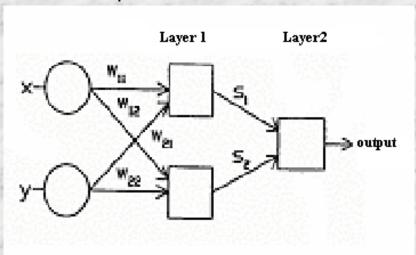
We spoke about the multilayer perceptron limitations. When the problem of linear separability was well understood it was also found that the problem of single-layer network can be overcome by adding more layers. For example, the two layer network may be formed by cascading two single-layer networks. These can perform more general classification, separating those points that are contained in convex open or closed regions.

A convex regions

A convex region – is one in which any two points in the region can be joined by a straight line that does not leave the region.



To understand the convexity limitation, consider a simple two-layer network with two inputs going to two neurons in the first layer, both feeding a single neuron in the layer 2. The output neuron threshold is set to 0.75 and weights s_i are both set to 0.5

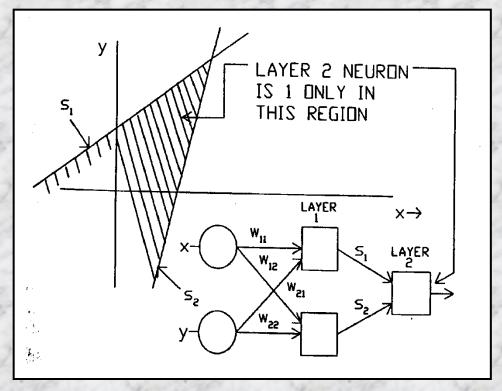


In this case the output of one (1) is required from both layer 1 neurons to exceed the threshold and to produce a one (1) on the output. Thus the output neuron performs a logical **AND** function.

Each neuron in layer 1 subdivides the **OXY** plane, producing tan output of one (1) on one side of the line.

The result of double subdivision when the output of one (1) of the layer 2 neuron is only over V-shaped

region.

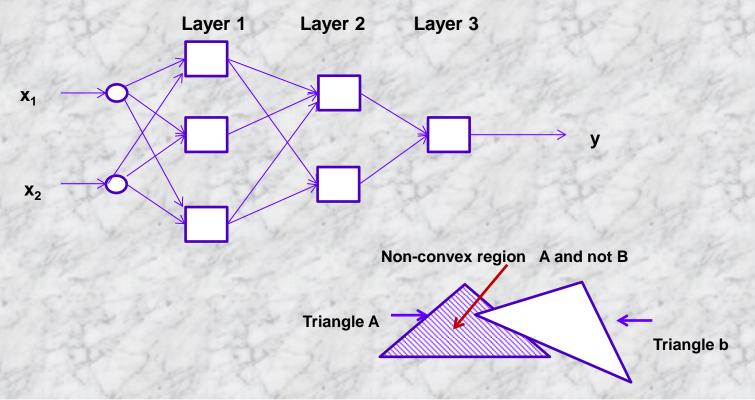


Similarly, three neurons used in the layer 1 further subdivide the plane, creating for example a triangle-shaped region. By including enough neurons in the layer 1, a convex region of any desired shape can be formed.

The layer 2 of course is not limited to the AND function, it can produce many other functions if the weights and threshold are suitably chosen.

A three-layer network is still more general, its limitation capability is limited only by the number of neurons and weights. There are no convexity constrains; the layer 3 neuron receives as input a group of convex polygons, and the logical combination of that need not to be convex.

A concave decision region formed by intersection of two convex regions. Logical operation A and not B



Existence (Kolmogorov) Theorem

Any continuous function of n variables can be computed using only linear summations and nonlinear but continuously increasing functions of only one variable. It effectively states that a three layer perceptron with n(2n+1) nodes using continuously increasing non-linearities can compute any continuous function of n variables. A three layer perceptron could then be used to create any continuous likelihood function required in a classifier.

Conclusion:

To create any arbitrary complex shape (decision region), we never need more that three layers in the network.

It gives the limitation on layers but does not define:

- how many elements is necessary to create a network (in general and in particular layers),
- how these elements should be connected,
- which weights value should be.

Network structure

Inconsistency in nomenclature

What is layer???

- some authors refer to the number of layers of variable weights
- some authors describe the number of layers of nodes

Usually, the nodes in the first layer, the input layer, merely distribute the inputs to subsequent layers, and do not perform and operations (summation or thresholding) n.b. some authors miss out these nodes.

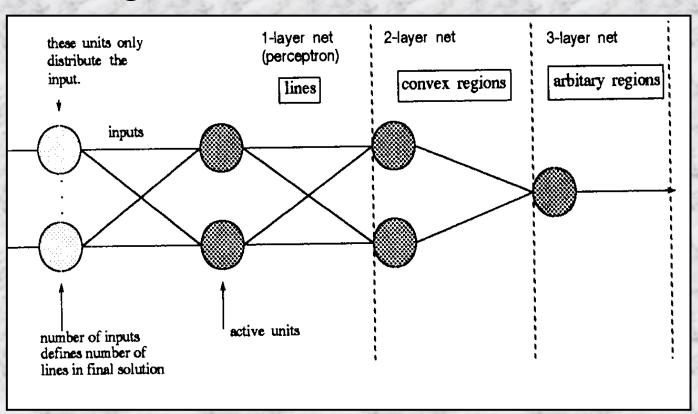
Network structure

What is a network layer?

A layer - it is the part of network structure which contains active elements performing some operation. A multilayer network receives a number of inputs. These are distributed by a layer of input nodes that do not perform any operation – these inputs are then passed along the first layer of adaptive weights to a layer of perceptron-like units, which do sum and threshold their inputs. This layer is able to produce classification lines in pattern space.

The output from this layer in then passed to another layer, and the output of this layer forms convex regions in pattern space. A further layer is able to define any arbitrary shape in pattern space.

Neural networks and their corresponding decision regions



Different non linearly separable problems and number of layers

Structure	Types of Decision Regions	Exclusive-OR Problem	Classes with Meshed regions	Most General Region Shapes
Single-Layer	Half Plane Bounded By Hyperplane	A B	B	
Two-Layer	Convex Open Or Closed Regions	A B A	B	
Three-Layer	Abitrary (Complexity Limited by No. of Nodes)	A B A	B	

Neural Networks - An Introduction Dr. Andrew Hunter

A multilayer perceptron is **fault-tolerant**, since it is distributed parallel processing element, with each node contributing to the final output.

If a node or its weights are lost or damaged. recall is impaired in quality, but the distributed nature of the information means that the damage has to be extensive before the network's response degrades badly. The network demonstrate graceful degradation rather that catastrophic failure.

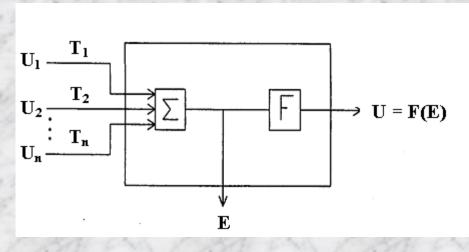


For many years there was no theoretical sound algorithm for training multilayer artificial neural networks. Since single-layer networks proved severely limited in what they could represent, the entire field went into virtual eclipse.

In 1986 Rumelhart and McClelland suggested a new learning rule known as *a backpropagation rule* which is today used in many practical applications like for example in solving the optimization problems.

The invention of the backpropagation algorithm has played a large part in the resurgence of interest in artificial neural networks.

Backpropagation is a systematic method for training multilayer artificial neural networks.

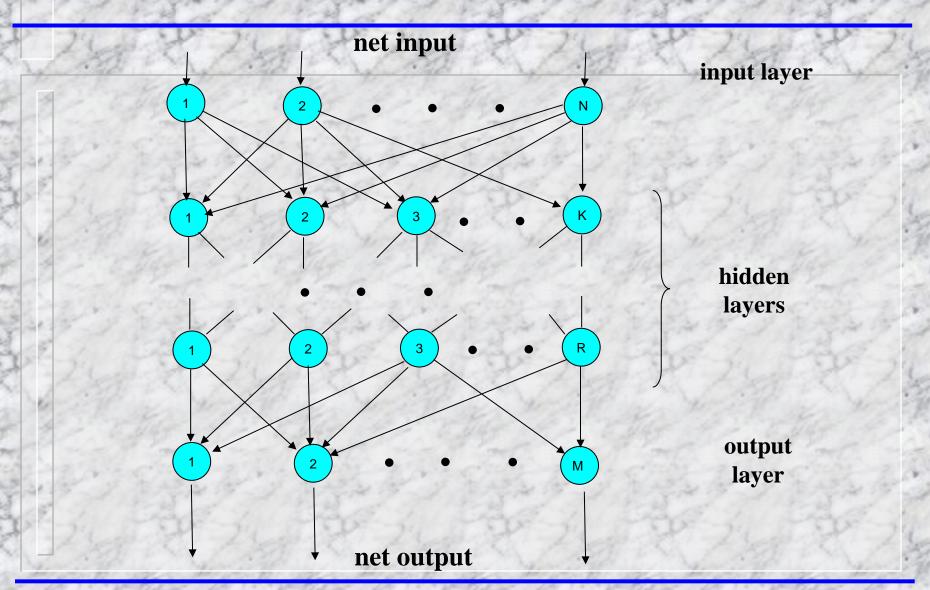


It is the rule how to change the weights T_{ij} between network elements.

The algorithm is based on the idea to minimize the square root of errors by use of the gradient descent method.

Assumptions:

- the net is the regular, multilayer structure
- the first layer input layer
- the last layer output layer
- layers between hidden layers
- feed forward propagation only



To define a state of j-th neuron in layer n we calculate the weighted sum of its \mathbf{M} inputs

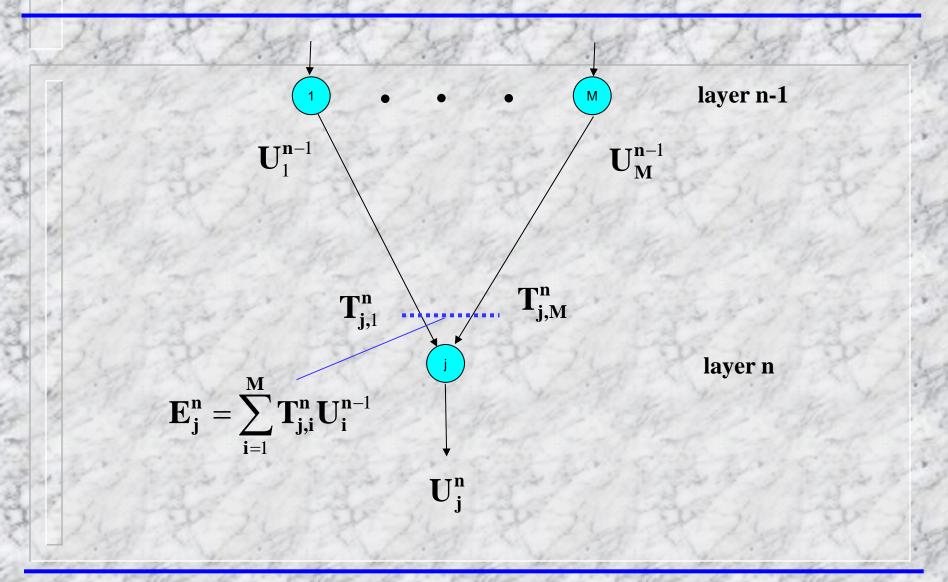
$$\mathbf{E}_{j}^{n} = \sum_{i=1}^{M} \mathbf{T}_{j,i}^{n} \mathbf{U}_{i}^{n-1} \tag{1}$$

where

 $\mathbf{E}_{\mathbf{i}}^{\mathbf{n}}$ weighted input sum of *j*-th neuron in layer *n*

 $\mathbf{T_{j,i}^n}$ connection weight between *i*-th neuron inlayer *n*-1 and *j*-th neuron in layer *n*

 $\mathbf{U}_{\mathbf{i}}^{\mathbf{n-1}}$ out put signal of *i*-th neuron in layer *n*-1



output signal of j-th neuron in layer n is defined by

$$\mathbf{U_j^n} = f(\mathbf{E_j^n}) \tag{2}$$

where f is the neuron transfer function

$$\mathbf{U_{j}^{n}} = f(\mathbf{E_{j}^{n}}) = \frac{1}{1 + exp[-(\mathbf{E_{j}^{n}} - \mathbf{\Theta_{j}^{n}})/\mathbf{\Theta_{0}}]}$$
(3)

f – sigmoid function, the output values \in (0;1),

 $\Theta^{\mathbf{n}}_{\mathbf{i}}$ - the threshold, $\,\Theta_{0}\,$ - slope

The global error **D** is a differentiable function of weights

$$D = \frac{1}{2} \cdot \sum_{j=1}^{M} (U_j^* - U_j^{out})^2$$
 (4)

where $\mathbf{U}_{\mathbf{j}}^{*}$ desired output (a target) of j-th neuron in the output layer

 U_{j}^{out} actual output of j-th neuron in the output layer

Aim: minimization of a global error **D**, modifying the network's weights.

Goal: define the rules of the weights adaptation to minimize the global error

Solution: the gradient descent method

$$\Delta \mathbf{T_{j,i}^n} = -\eta \frac{\partial \mathbf{D}}{\partial \mathbf{T_{j,i}^n}}$$
 (5)

where η is the learning coefficient

Each weight is changed according to the value and direction of negative gradient on the hyperplane **D(T)**.

The partial derivative in (5) can be calculated by use of the chain rule.

$$\frac{\partial \mathbf{D}}{\partial \mathbf{T_{j,i}^n}} = \frac{\partial \mathbf{D}}{\partial \mathbf{E_j^n}} \cdot \frac{\partial \mathbf{E_j^n}}{\partial \mathbf{T_{j,i}^n}} \tag{6}$$

We define the change in error as a function of the change in the net inputs to a *j*-th neuron as

$$\mathbf{d_{j}^{n}} = -\frac{\partial \mathbf{D}}{\partial \mathbf{E_{j}^{n}}}$$

$$\frac{\partial E_{j}^{n}}{\partial T_{j,i}^{n}} = -\frac{\partial}{\partial T_{j,i}^{n}} \sum_{k=1}^{M} T_{j,k}^{n} U_{k}^{n-1} = \sum_{k=1}^{M} \frac{\partial T_{j,k}^{n}}{\partial T_{j,i}^{n}} U_{k}^{n-1} = U_{i}^{n-1}$$
(7)

because

$$\frac{\partial T^n_{j,k}}{\partial T^n_{j,i}} = \begin{cases} 0 \text{ for } k \neq i \\ 1 \text{ for } k = i \end{cases}$$

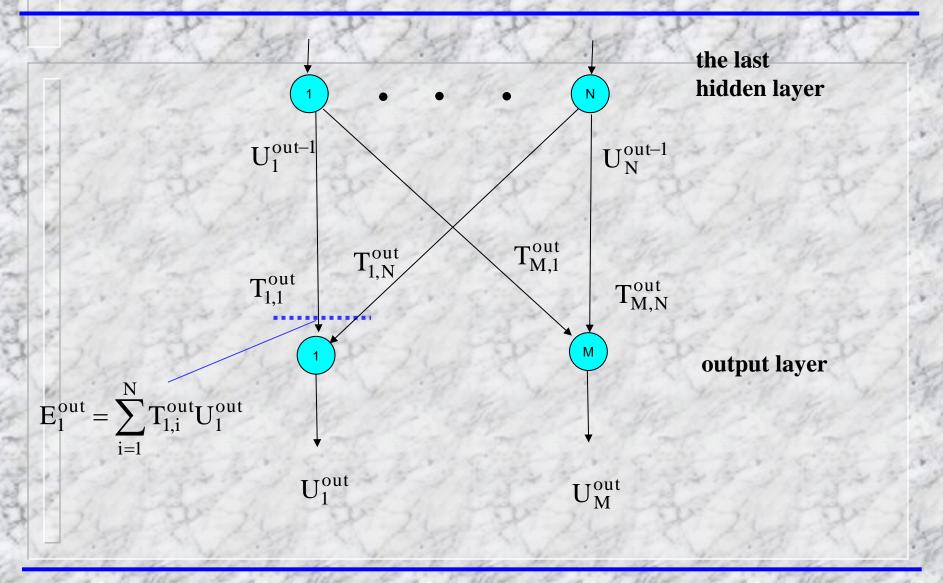
Finally, we get

$$\Delta \mathbf{T}_{j,i}^{n} = \boldsymbol{\eta} \mathbf{d}_{j}^{n} \mathbf{U}_{i}^{n-1}$$
 (8)

For the elements in the output layer

$$\Delta \mathbf{T_{j,i}^{out}} = \eta \mathbf{d_{j}^{out}} \mathbf{U_{i}^{out-1}}$$
 (9)

Decrease of **D** means the changes proportional to $\mathbf{d}_{\mathbf{j}}^{\text{out}}\mathbf{U}_{\mathbf{i}}^{\text{out-1}}$



Now, we need to know what d_{j}^{out} is for each of the units

$$\mathbf{d_{j}^{out}} = -\frac{\partial \mathbf{D}}{\partial \mathbf{E_{j}^{out}}} = \frac{\partial \mathbf{D}}{\partial \mathbf{U_{j}^{out}}} \frac{\partial \mathbf{U_{j}^{out}}}{\partial \mathbf{E_{j}^{out}}}$$
(10)

where

$$\frac{\partial \mathbf{U_{j}^{out}}}{\partial \mathbf{E_{j}^{out}}} = \mathbf{f'}(\mathbf{E_{j}^{out}})$$

differentiate D with respect to $\mathbf{U}_{\mathbf{j}}^{out}$ giving

$$\frac{\partial \mathbf{D}}{\partial \mathbf{U_{j}^{out}}} = -(\mathbf{U_{j}^{*}} - \mathbf{U_{j}^{out}})$$
 (11)

thus

$$\mathbf{d_j^{out}} = (\mathbf{U_j^*} - \mathbf{U_j^{out}}) \cdot \mathbf{f'}(\mathbf{E_j^{out}})$$
 (12)

This is useful for the output units to modify weights $T_{j,i}^{out}$ between the neurons of the output layer and the neurons of the last hidden layer (since the target and output both are available)

$$T_{j,i}^{out}(t) = T_{j,i}^{out}(t-1) + \tau \Delta T_{j,i}^{out}$$

where τ is a learning coefficient (13)

The formula (9) can be rewritten to avoid oscillations

$$\Delta T_{j,i}^{\text{out}}(t) = \alpha \Delta T_{j,i}^{\text{out}}(t-1) \qquad (1-\alpha)d_j^{\text{out}}f(E_i^{\text{out-1}}) \qquad (14)$$

where $\mathbf{E}_{\mathbf{j}}^{\mathrm{wy-1}}$ is the weighted input sum of *i*-th element from the last hidden layer, α (so called smoothing parameter), is a constant value defining the effect of the previous weights modification on the actual modification .

This method is very useful for the output layer elements because we have access both to the output signal, target signal.

However, for the elements located in the hidden layers unfortunately does not work.

If *j*-th neuron **does not belong** to the output layer but is an element of a hidden layer *s*, then its local weight modification is calculated from

$$\mathbf{d}_{\mathbf{j}}^{\mathbf{s}} = -\frac{\partial \mathbf{D}}{\partial \mathbf{E}_{\mathbf{j}}^{\mathbf{s}}}$$
 (15)

Because

$$\frac{\partial \mathbf{D}}{\partial \mathbf{E_{j}^{n}}} = \frac{\partial \mathbf{D}}{\partial \mathbf{U_{j}^{n}}} \cdot \frac{\partial \mathbf{U_{j}^{n}}}{\partial \mathbf{E_{j}^{n}}} = f'(\mathbf{E_{j}^{n}}) \frac{\partial \mathbf{D}}{\partial \mathbf{U_{j}^{n}}}$$
(16)

also

$$\left| \frac{\partial \mathbf{D}}{\partial \mathbf{U_j^n}} = \sum_{k=1}^{N} \frac{\partial \mathbf{D}}{\partial \mathbf{E_k^{n+1}}} \frac{\partial \mathbf{E_k^{n+1}}}{\partial \mathbf{U_j^n}} = \sum_{k=1}^{N} \left[\frac{\partial \mathbf{D}}{\partial \mathbf{E_k^{n+1}}} \frac{\partial}{\partial \mathbf{U_j^n}} \left(\sum_{i=1}^{M} \mathbf{T_{k,i}^{n+1}} \mathbf{U_i^n} \right) \right] = \sum_{k=1}^{N} \left[\frac{\partial \mathbf{D}}{\partial \mathbf{E_k^{n+1}}} \frac{\partial}{\partial \mathbf{U_j^n}} \left(\sum_{i=1}^{M} \mathbf{T_{k,i}^{n+1}} \mathbf{U_i^n} \right) \right] = \sum_{k=1}^{N} \left[\frac{\partial \mathbf{D}}{\partial \mathbf{E_k^{n+1}}} \frac{\partial}{\partial \mathbf{E_k^{n+1}}} \frac{\partial}{\partial \mathbf{E_k^{n+1}}} \left(\sum_{i=1}^{M} \mathbf{T_{k,i}^{n+1}} \mathbf{U_i^n} \right) \right] = \sum_{k=1}^{N} \left[\frac{\partial \mathbf{D}}{\partial \mathbf{E_k^{n+1}}} \frac{\partial}{\partial \mathbf{E_k^{n+1}}} \frac{\partial}{\partial \mathbf{E_k^{n+1}}} \left(\sum_{i=1}^{M} \mathbf{T_{k,i}^{n+1}} \mathbf{U_i^n} \right) \right] = \sum_{k=1}^{N} \left[\frac{\partial \mathbf{D}}{\partial \mathbf{E_k^{n+1}}} \frac{\partial}{\partial \mathbf{E_k^{n+1}}} \frac{\partial}{\partial \mathbf{E_k^{n+1}}} \frac{\partial}{\partial \mathbf{E_k^{n+1}}} \frac{\partial}{\partial \mathbf{E_k^{n+1}}} \frac{\partial}{\partial \mathbf{E_k^{n+1}}} \right] = \sum_{k=1}^{N} \left[\frac{\partial \mathbf{D}}{\partial \mathbf{E_k^{n+1}}} \frac{\partial}{\partial \mathbf{E_k^{n+1$$

$$=\sum_{k=1}^{N}\frac{\partial \mathbf{D}}{\partial \mathbf{E}_{\mathbf{k}}^{\mathbf{n+1}}}\mathbf{T}_{\mathbf{k},\mathbf{j}}^{\mathbf{n+1}}=\sum_{k=1}^{N}\mathbf{d}_{\mathbf{k}}^{\mathbf{n+1}}\mathbf{T}_{\mathbf{k},\mathbf{j}}^{\mathbf{n+1}}$$
(17)

$$\mathbf{d_{j}^{n}} = -\frac{\partial \mathbf{D}}{\partial \mathbf{E_{j}^{n}}} = f'(\mathbf{E_{j}^{n}}) \sum_{k=1}^{N} \mathbf{d_{j}^{n+1}} \mathbf{T_{k,j}^{n+1}}$$
(18)

it modify the weight $T_{j,i}^n$ between *i*-th element of the layer n-1 and *j*-th element of the layer n

$$T_{j,i}^{n}(t) = T_{j,i}^{n}(t-1) + \tau \Delta T_{j,i}^{n}$$
 (19)

a

$$\Delta \mathbf{T}_{\mathbf{j},\mathbf{i}}^{\mathbf{n}}(\mathbf{t}) = \alpha \Delta \mathbf{T}_{\mathbf{j},\mathbf{i}}^{\mathbf{n}}(\mathbf{t} - \mathbf{1}) + (1 - \alpha) \, \mathbf{d}_{\mathbf{j}}^{\mathbf{n}} f(\mathbf{E}_{\mathbf{i}}^{\mathbf{n}-1})$$
(20)

Iterative calculations correct weight "backwards", to the input layer.

This learning procedure is repeated for each learning pattern, until the network will generate the correct answer for every input signal (pattern) – of course with the defined accuracy.

Example:

If
$$f(\mathbf{E}) = \frac{1}{1 + exp(-k\mathbf{E})}$$

$$f(E) \in (0;1), k>0$$

when $k \Rightarrow \infty$
then $f(E) \Rightarrow$ step function

$$\mathbf{U_{j}^{out}} = f(\mathbf{E_{j}^{out}}) = \frac{1}{1 + \exp(-k\mathbf{E_{j}^{out}})}$$

calculating the derivative f'(E) for j-th element:

$$f'(\mathbf{E}) = \frac{ke^{-k\mathbf{E}}}{(1+e^{-k\mathbf{E}})^2} = kf(\mathbf{E}) \cdot [1-f(\mathbf{E})] =$$
$$= k\mathbf{U}_{\mathbf{j}}^{\mathbf{out}} (1-\mathbf{U}_{\mathbf{j}}^{\mathbf{out}})$$

(this derivative simplify he calculations)

for the elements in the output layer

$$\begin{aligned} \mathbf{d_j^{out}} &= (\mathbf{U_j^*} - \mathbf{U_j^{out}}) \cdot \mathbf{f'(E_j^{out})} = \\ &= \mathbf{k} \cdot \mathbf{U_j^{out}} \cdot (1 - \mathbf{U_j^{out}}) \cdot (\mathbf{U_j^*} - \mathbf{U_j^{out}}) \end{aligned}$$

for the elements in the hidden layers

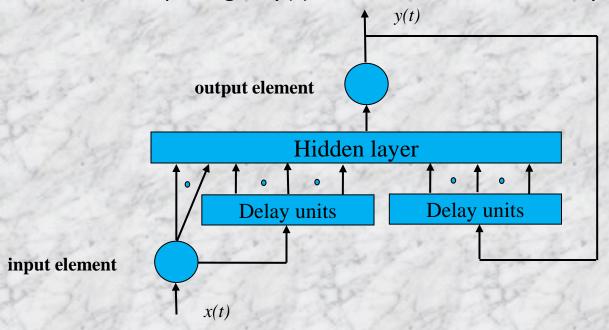
$$\mathbf{d}_{j}^{n} = k \cdot \mathbf{U}_{j}^{n} \cdot (1 - \mathbf{U}_{j}^{n}) \cdot \sum_{k=1}^{N} \mathbf{d}_{k}^{n+1} \mathbf{T}_{j,k}^{n+1}$$



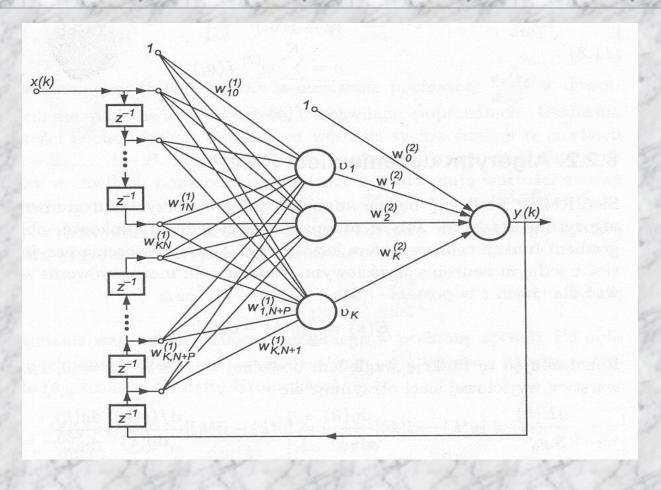
Recurrent neural network architectures consists of a standard Multi-Layer Perceptron (MLP) plus added loops.

<u>Assumption</u>: only one input nod (input signal x(t)) and one output nod (output signal y(t)), and one hidden layer.

The network input signal is composed from input x, delayed inputs and external recurrence output signal y(t) with some unit time delays.



Structure of the RMLP network



RMLP is a dynamic network with delayed input and output signals.

The farther analysis is performed for only one input node (signal x(t)), one output node (signal y(t)) and one hidden layer.

The function

$$y(t+1)=f(x(t),x(t-1),...,x(t-(N-1)),y(t-1),y(t-2),...,y(t-P))$$

where

N-1 \Rightarrow number of delayed input signals,

P ⇒ number of delayed output signals,

 $K \Rightarrow$ number of neurons in the hidden layer,

Vector input signal x applied to the network input in the time t

$$x(t)=[1,x(t),x(t-1),...,x(t-(N-1)),y(t-P),y(t-P+1),...,y(t-1)]$$

Denote

$$u_i = \sum_{j=0}^{N+P} w_{ij}^{(1)} x_j$$

$$v_i = f(u_i)$$

$$g = \sum_{i=0}^{K} w_i^{(2)} f(u_j)$$

$$y = f(g)$$

weighted input signal of the i-th neuron of the hidden layer

input signal of the *i*-th neuron of the hidden layer

weighted input signal of the output neuron

output signal of the output neuron

(1-4)



RMLP Learning Algorythm

The gradient learning algorithm is used. The gradient of ab error function with respect to each networks' weight is calculated. For RMLP network the error function can be define by:

$$E(t) = \frac{1}{2} [y(t) - d(t)]^{2}$$

Differentiating with respect to each weight $w_{\alpha}^{(2)}$ in the second layer

$$\frac{\partial E(t)}{\partial w_{\alpha}^{(2)}} = [y(t) - d(t)] \frac{dy(t)}{dw_{\alpha}^{(2)}} = [y(t) - d(t)] \frac{df(g(t))}{dg(t)} \frac{dg(t)}{dw_{\alpha}^{(2)}} =$$

$$= [y(t) - d(t)] \frac{df(g(t))}{dg(t)} \sum_{i=0}^{K} \frac{d(w_{\alpha}^{(2)}v_{i}(t))}{dw_{\alpha}^{(2)}} =$$

$$= [y(t) - d(t)] \frac{df(g(t))}{dg(t)} \left[v_{\alpha}(t) + \sum_{i=0}^{K} w_{\alpha}^{(2)} \frac{dv_{i}(t)}{dw_{\alpha}^{(2)}} \right]$$

because

$$\frac{d(w_{\alpha}^{(2)})}{dw_{\alpha}^{(2)}} = \begin{cases} 1 & \text{for } i = \alpha \\ 0 & \text{for } i \neq \alpha \end{cases}$$

$$\begin{split} &\frac{dv_{i}(t)}{dw_{\alpha}^{(2)}} = \frac{df\left(u_{i}(t)\right)}{du_{i}(t)} \sum_{j=0}^{N+P} w_{ij}^{(1)} \frac{dx_{j}}{dw_{\alpha}^{(2)}} = \\ &= \frac{df\left(u_{i}(t)\right)}{du_{i}(t)} \sum_{j=N+1}^{N+P} w_{ij}^{(1)} \frac{dy(t-P-1+(j-N))}{dw_{\alpha}^{(2)}} = \\ &= \frac{df\left(u_{i}(t)\right)}{du_{i}(t)} \sum_{j=0}^{P} w_{ij}^{(1)} \frac{dy(t-P-1+j)}{dw_{\alpha}^{(2)}} \end{split}$$

next

$$\frac{dy(t)}{dw_{\alpha}^{(2)}} =$$

$$= \frac{df(g(t))}{dg(t)} \left[v_{\alpha}(t) + \sum_{i=0}^{K} w_{i}^{(2)} \frac{df(u_{i}(t))}{du_{i}(t)} \sum_{j=1}^{P} w_{i,j+N}^{(1)} \frac{dy(t-P-1+j)}{dw_{\alpha}^{(2)}} \right]$$

this formula enable calculater the $\frac{dy(t)}{dw_{\alpha}^{(2)}}$ in the epoch t on the base of its previus values

Using the gradient descent method the weight change of the output layer is defined by

$$\Delta w_{\alpha}^{(2)} = -\eta [y(t) - d(t)] \frac{dy(t)}{dw_{\alpha}^{(2)}}$$
 (5)

Similarly can be calculated the weight change in the hidden layer. $w^{(1)}_{\alpha,\beta}$

After calualtion of the derivative of a signal y(t) with respect of the weight ijn the hidden layer

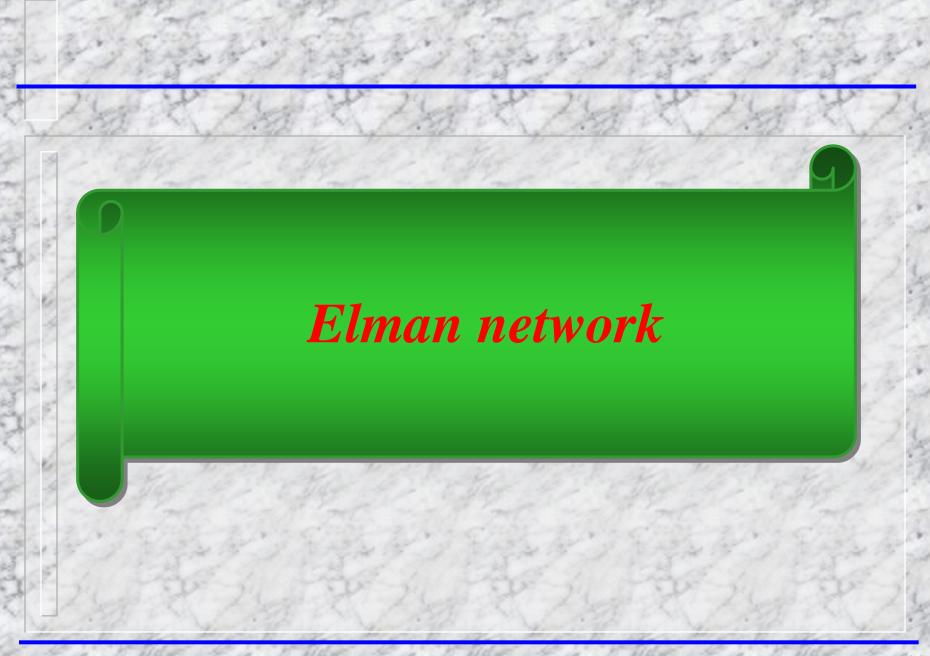
$$\frac{dy(t)}{dw_{\alpha,b}^{(1)}} = \frac{df(g(t))}{dg(t)} \sum_{i=0}^{K} w_i^{(2)} \frac{df(u_i(t))}{du_i(t)} \left[\sum_{j=1}^{P} w_{i,j+N}^{(1)} \frac{dy(t-P-1+j)}{dw_{\alpha}^{(2)}} \right]$$

and the formula for the weight $w_{\alpha,\beta}^{(1)}$ change in the hidden layer

$$\Delta w_{\alpha,\beta}^{(1)} = -\eta [y(t) - d(t)] \frac{dy(t)}{dw_{\alpha,\beta}^{(1)}}$$
 (6)

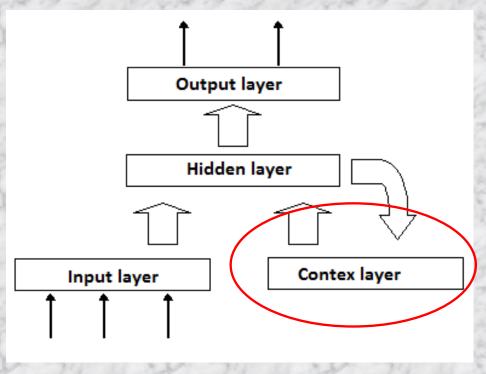
Summary of a learning algorithm

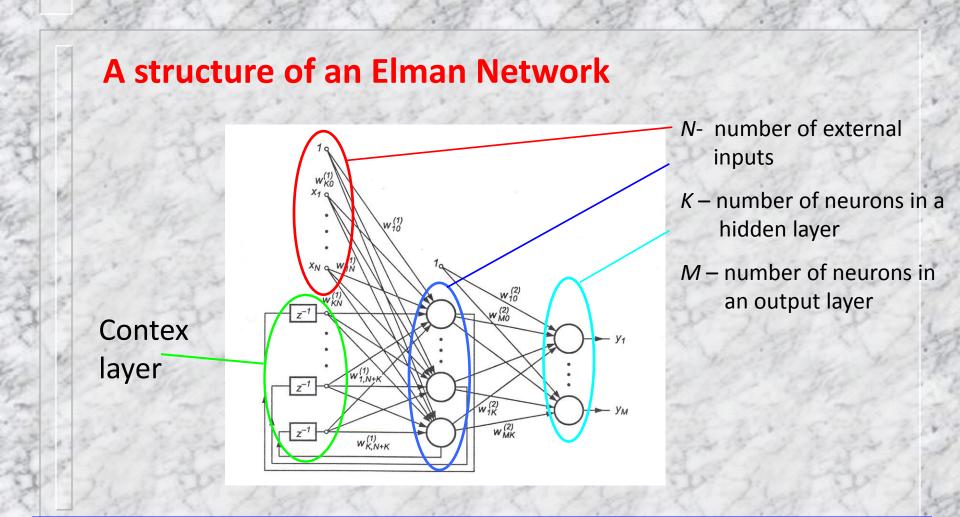
- Initialize the weight vectors in the hidden layer and output layer
- 2. Evaluate the neuron state in epoch (t) with input x (1-4)
 - 3. Calculates the values of $\frac{dy(t)}{dw_{\alpha}^{(2)}}$ and $\frac{dy(t)}{dw_{\alpha\beta}^{(1)}}$ for every $\alpha\beta$
 - 4. Updates the weights according to 5 i 6
 - 5. Go to step 2 algorithm.



An Elman network is an MLP with a single hidden layer and in addition it contains connections from the hidden layer's neurons to the context units. The context units store the output values from the hidden neurons in a time unit and these values are fed as additional inputs to the hidden neurons in the next time unit.

Elman Network – SRN (Simple Recurrent Network) is a simplification of Multi Layer Perceptron





Vector input signal:

$$[1, x_1^{(1)}(t), ..., x_N^{(1)}(t), x_{N+1}^{(1)}(t), ..., x_{N+K}^{(1)}(t)]$$

the elements denoted j=N+1,...,N+K are from the hidden layer from the previous epoch

$$[1, x_1^{(1)}(t), ..., x_N^{(1)}(t), y_1^{(1)}(t-1), ..., y_K^{(1)}(t-1)]$$

Notation

$$u_i(t) = \sum_{j=0}^{N+K} w_{ij}^{(1)} x_j(t)$$

$$v_i = f_1(u_i(t))$$

$$g_i(t) = \sum_{j=0}^K w_{ij}^{(2)} v_j(t)$$

$$y_i = f_2(g_i(t))$$

weighted input signal of the *i*-th neuron in a hidden layer

output signal of the *i*-th neuron of a hidden layer

weighted input signal of the *i*-th neuron of the output layer

output signal of the *i*-th neuron of the output layer

Elman learning algorithm

The network will be learn according to steepest descent algorithm. Similarly to feedforward network a gradient of the cost function will be calculated with respect to every network's weight. For the Elman network a cost function can be defined by

$$E(t) = \frac{1}{2} \sum_{i=1}^{M} [y_i(t) - d_i(t)]^2 = \frac{1}{2} \sum_{i=1}^{M} [e_i(t)]^2$$

Differentiating this function with respect to any output layer weight $w_{\alpha\beta}^{(2)}$ we obtain

$$\begin{split} &\nabla^{(2)}_{\alpha\beta}E(t) = \frac{\partial E(t)}{\partial w^{(2)}_{\alpha\beta}} = \sum_{i=1}^{M} [y_i(t) - d_i(t)] \frac{df_2(g_i(t))}{dg_i(t)} \frac{dg_i(t)}{dw^{(2)}_{\alpha\beta}} = \\ &= \sum_{i=1}^{M} [y_i(t) - d_i(t)] \frac{df_2(g_i(t))}{dg_i(t)} \sum_{j=0}^{K} \frac{d(w^{(2)}_{ij}v_j(t))}{dw^{(2)}_{\alpha\beta}} = \\ &= \sum_{i=1}^{M} [y_i(t) - d_i(t)] \frac{df_2(g_i(t))}{dg_i(t)} \sum_{j=0}^{K} \left[\frac{d(v_j(t))}{dw^{(2)}_{\alpha\beta}} w^{(2)}_{ij} + \frac{d(w^{(2)}_{ij}(t))}{dw^{(2)}_{\alpha\beta}} v_j(t) \right] \end{split}$$

Because connections between the hidden layer and output layer are unidirectional

$$\frac{d(v_j(t))}{dw_{\alpha\beta}^{(2)}} = 0$$

yields to

$$\nabla_{\alpha\beta}^{(2)} E(t) = \frac{\partial E(t)}{\partial w_{\alpha\beta}^{(2)}} = \sum_{i=1}^{M} [y_i(t) - d_i(t)] \frac{df_2(g_i(t))}{dg_i(t)} \sum_{j=0}^{K} \frac{d(w_{ij}^{(2)}(t))}{dw_{\alpha\beta}^{(2)}} v_j(t)$$
(7)

According to the method of steepest descent the weights' change in output layer are defined by

$$w_{\alpha\beta}^{(2)}(t+1) = w_{\alpha\beta}^{(2)}(t) - \eta \nabla_{\alpha\beta}^{(2)} E(t)$$
(8)

Weights change in the hidden layer is more complicated because of the feedbacks.

$$\nabla_{\alpha\beta}^{(1)} E(t) = \frac{\partial E(t)}{\partial w_{\alpha\beta}^{(1)}} = \sum_{i=1}^{M} [y_i(t) - d_i(t)] \frac{df_2(g_i(t))}{dg_i(t)} \sum_{j=0}^{K} \frac{d(w_{ij}^{(2)} v_j(t))}{dw_{\alpha\beta}^{(1)}} = \sum_{i=1}^{M} [y_i(t) - d_i(t)] \frac{df_2(g_i(t))}{dg_i(t)} = \sum_{j=0}^{M} \frac{d(w_{ij}^{(2)} v_j(t))}{dw_{\alpha\beta}^{(1)}} = \sum_{i=1}^{M} [y_i(t) - d_i(t)] \frac{df_2(g_i(t))}{dg_i(t)} = \sum_{j=0}^{M} \frac{d(w_{ij}^{(2)} v_j(t))}{dw_{\alpha\beta}^{(1)}} = \sum_{j=0}^{M} \frac{d(w_{ij}^{(2)} v_j(t))}{dw_{\alpha\beta}^{(2)}} = \sum_{j=0}^{M} \frac{d(w_{ij}^{(2)} v_j(t)}{dw_{\alpha\beta}^{(2)}} = \sum_{j=0}^{M} \frac{d(w_{ij}^{(2)} v_j(t)}{dw_{\alpha\beta}^{(2)}} = \sum_{j=0}^{M} \frac{$$

$$= \sum_{i=1}^{M} [y_i(t) - d_i(t)] \frac{df_2(g_i(t))}{dg_i(t)} \sum_{j=0}^{K} \left[\frac{d(v_j(t))}{dw_{\alpha\beta}^{(1)}} w_{ij}^{(2)} \right]$$
(9)

$$\frac{dv_{j}(t)}{dw_{\alpha\beta}^{(1)}} = \frac{df_{1}(u_{j})}{du_{j}} \sum_{m=0}^{N+K} \frac{d(x_{m}(t)w_{jm}^{(1)})}{dw_{\alpha\beta}^{(1)}} = \frac{df_{1}(u_{j})}{du_{j}} \left[\mathcal{S}_{j\alpha}x_{\beta} + \sum_{m=0}^{N+K} \frac{dx_{m}(t)}{dw_{\alpha\beta}^{(1)}} w_{jm}^{(1)} \right]$$

Elman Network

and next

$$\frac{dv_{j}(t)}{dw_{\alpha\beta}^{(1)}} = \frac{df_{1}(u_{j})}{du_{j}} \left[\delta_{j\alpha} x_{\beta} + \sum_{m=N+1}^{N+K} \frac{dv_{(m-N)}(t-1)}{dw_{\alpha\beta}^{(1)}} w_{jm}^{(1)} \right] = \frac{df_{1}(u_{j})}{du_{j}} \left[\delta_{j\alpha} x_{\beta} + \sum_{m=1}^{K} \frac{dv_{m}(t-1)}{dw_{\alpha\beta}^{(1)}} w_{j,m+N}^{(1)} \right]$$
(10)

The last formula (10) allows to calculate the derivatives of cost function with respect to weights of hidden layer in moment *t*. It is the recurrent formula defining derivative in a moment *t* dependence of a derivative in a moment *t*-1.

Elman Network

Assuming the initial values in a moment t = 0

$$\frac{dv_1(0)}{dw_{\alpha\beta}^{(2)}} = \frac{dv_2(0)}{dw_{\alpha\beta}^{(2)}} = \dots = \frac{dv_K(0)}{dw_{\alpha\beta}^{(2)}} = 0$$

and basing on the steepest descent method the weights' change in the hidden layer is defined by

$$w_{\alpha\beta}^{(1)}(t+1) = w_{\alpha\beta}^{(1)}(t) - \eta \nabla_{\alpha\beta}^{(1)} E(t)$$
(11)

Elman Network

Elman training algorithm

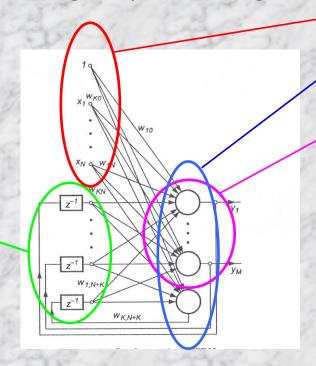
- Initialize the weight vector with random values in , the learning rate , the repetitions counter () and the epochs counter (). Initialize the context nodes at 0.5.
- 2.Let the network's weight vector in the beginning of epoch
 - 2.1 Start of epoch. Store the current values of the weight vector
 - 2.2 For n=1,2,...,N
 - 2.2.1 Select the training example (xⁿ, tⁿ)and apply the error backpropagation in order to compute the partial derivatives $\eta \frac{\partial E^n}{\partial w_i}$
 - 2.2.2 Update the weights $w_i(k+1) = w_i(k) \eta \frac{\partial E^n}{\partial w_i}$
 - 2.2.3 Copy the hidden nodes' values to the context units.
 - 2.3 End of epoch k. Termination check. If true, terminate.
- 3. k=k+1. Go to step 2.



RTRN Network (Real Time Recurrent Network) is the network for real-time signal processing

Network structure

K element context layer



N- number of inputs
K – number of neurons in hidden layer
M – number of output neurons (from the K hidden elements)

$$u_i(t) = \sum_{j=0}^{N+K} w_{ij}^{(1)} x_j(t)$$

weighted input signal of i-th neuron in hidden layer

$$y_i = f(u_i(t))$$

output signal from *i*-th neuron in hidden layer (13)

Vector input signal x(t) and delayed one cycle vector y(t-1) create the excitation network signal

$$[1, x_1(t), ..., x_2(t), x_N(t), y_1(t-1), ..., y_K(t-1)]$$
(14)

This is the special case of the Elman network with constant values of output weights $(w_{\alpha}^{(2)})$ and defined by

$$w_{ij}^{(2)} = \delta_{ij} = \begin{cases} 1 & \text{for } i = j \\ 0 & \text{for } i \neq j \end{cases}$$

Modifying the Elman's algorithms one obtain the new learning algorithm for that network

RTRN training algorithm

- 1. Initialization of weights (usually uniform distribution from [-1;1]).
- 2. Calculation of naurons' state in (t=0,1,...) (signals u_i and y_i formulas (12) and (13) and next vector excitation signal (14).
- 3. Calculation of $\frac{dy_j(t)}{dw_{\alpha\beta}} = 0$ according to

$$\frac{dy_{j}(t)}{dw_{\alpha\beta}} = \frac{df(u_{j})}{du_{j}} \left[\delta_{j\alpha} + \sum_{t=1}^{K} \frac{dy_{j}(t-1)}{dw_{\alpha\beta}} w_{j,t+N} \right]$$

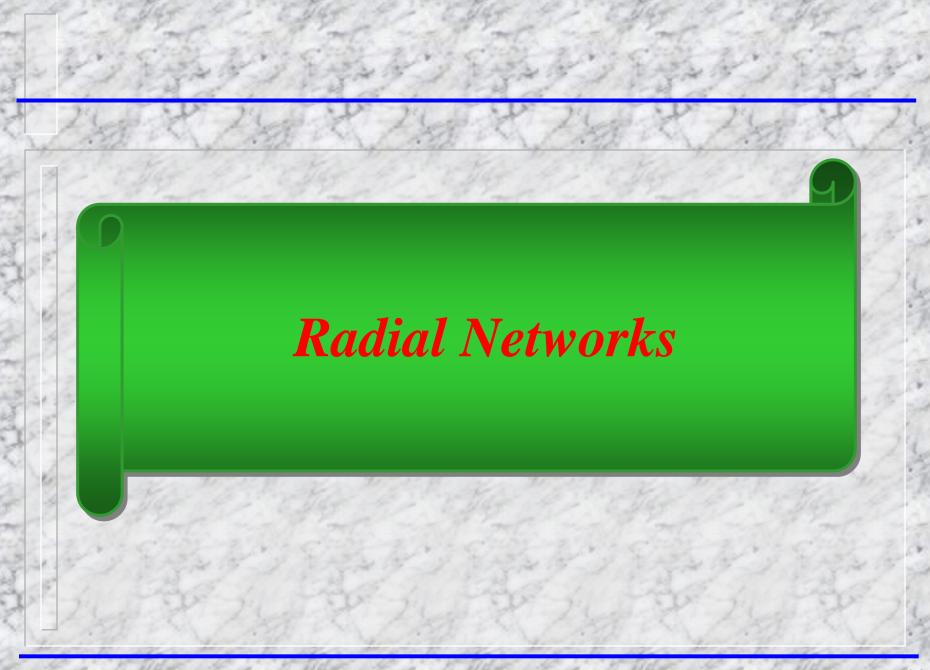
RTRN training algorithm

4. Upgrades of weights according to steepest descent algorithm according to

$$w_{\alpha\beta}(t+1) = w_{\alpha\beta}(t) - \eta \sum_{j=1}^{M} [y_j(t) - d_j(t)] \frac{dy_j(t)}{dw_{\alpha\beta}}$$

for
$$\alpha = 1, 2, ..., K$$
, $\beta = 0, 1, 2,, N+K$

5. Go to step 2.



Network structure

Radial networks

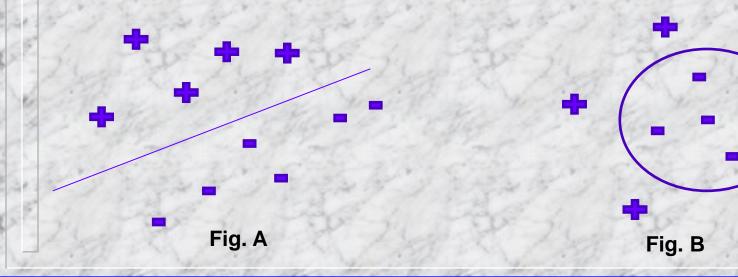
For special purposes sometimes we are using the radial neurons i.e. neurons with *Radial Basis Function RBF*. They have non typical aggregation methods of input data, they uses non typical transfer function (Gauss function) and they learned by the special way.

Data aggregation consist on the calculation of a distance between an input signal (vector **X**), and established in a learning process centroid of a certain set **T**.

Network structure

Sigmoid and radial neuron

Sigmoidal neuron represents in the multidimensional space hyperplane separating space into two categories (Fig.A). Radial neuron represents hypersphere performing circular separation around the central point (Fig.B).



Radial basis function network

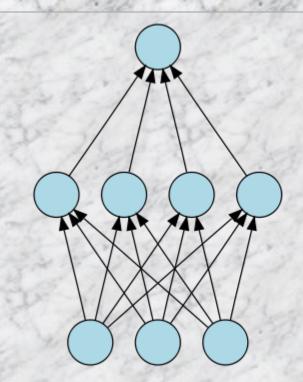
Radial basis function network (RBF) uses radial basis functions as activation functions typically have three layers: an input layer, a hidden layer with a non-linear RBF activation function and a linear output layer.

Functions that depend only on the distance from a center vector are radially symmetric about that vector, hence the name *radial basis function*. In the basic form all inputs are connected to each hidden neuron. The norm is typically taken to be the Euclidean distance and the radial basis function is commonly taken to be Gaussian.

Radial basis function network

The output of the network is then a scalar function of the input vector, and is given $\mathbf{b}(\mathbf{x}) = \sum_{i=1}^{N} a_i \rho(||\mathbf{x} - \mathbf{c}_i||)$

where N is the number of neurons in the hidden layer, C_i is the center vector for neuron i and a_i is the weight of neuron i in the linear output neuron. Functions that depend only on the distance from a center vector are radially symmetric about that vector. The radial basis function is commonly taken to be Gaussian.



Output y

Linear weights

Radial basis functions

Weights

Input x

$$\rho(\|\mathbf{x} - \mathbf{c}_i\|) = \exp[-\beta \|\mathbf{x} - \mathbf{c}_i\|^2]$$