

ASCI 2000

Learned from Neural Networks

R.P.W. Duin

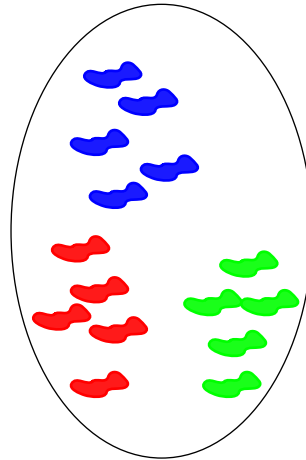
*Pattern Recognition Group
Delft University of Technology
The Netherlands*

duin@tn.tudelft.nl

Lommel, June 2000

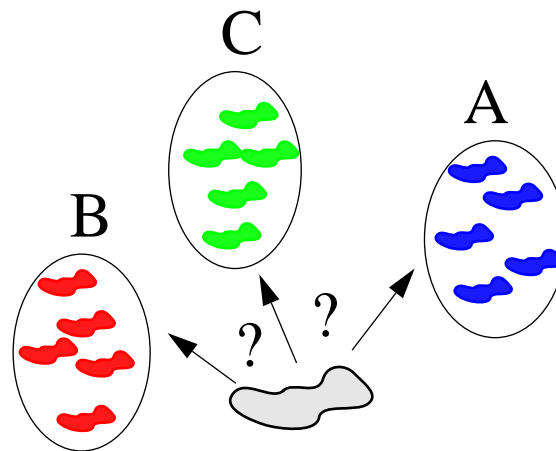
The Pattern Recognition Problem

Unsupervised Problem
(Clustering)



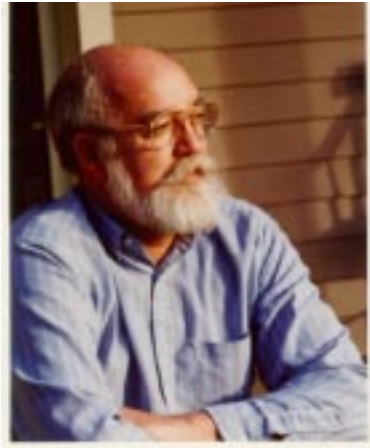
Define the classes
for a set of real world objects

Supervised Problem
(Classification)



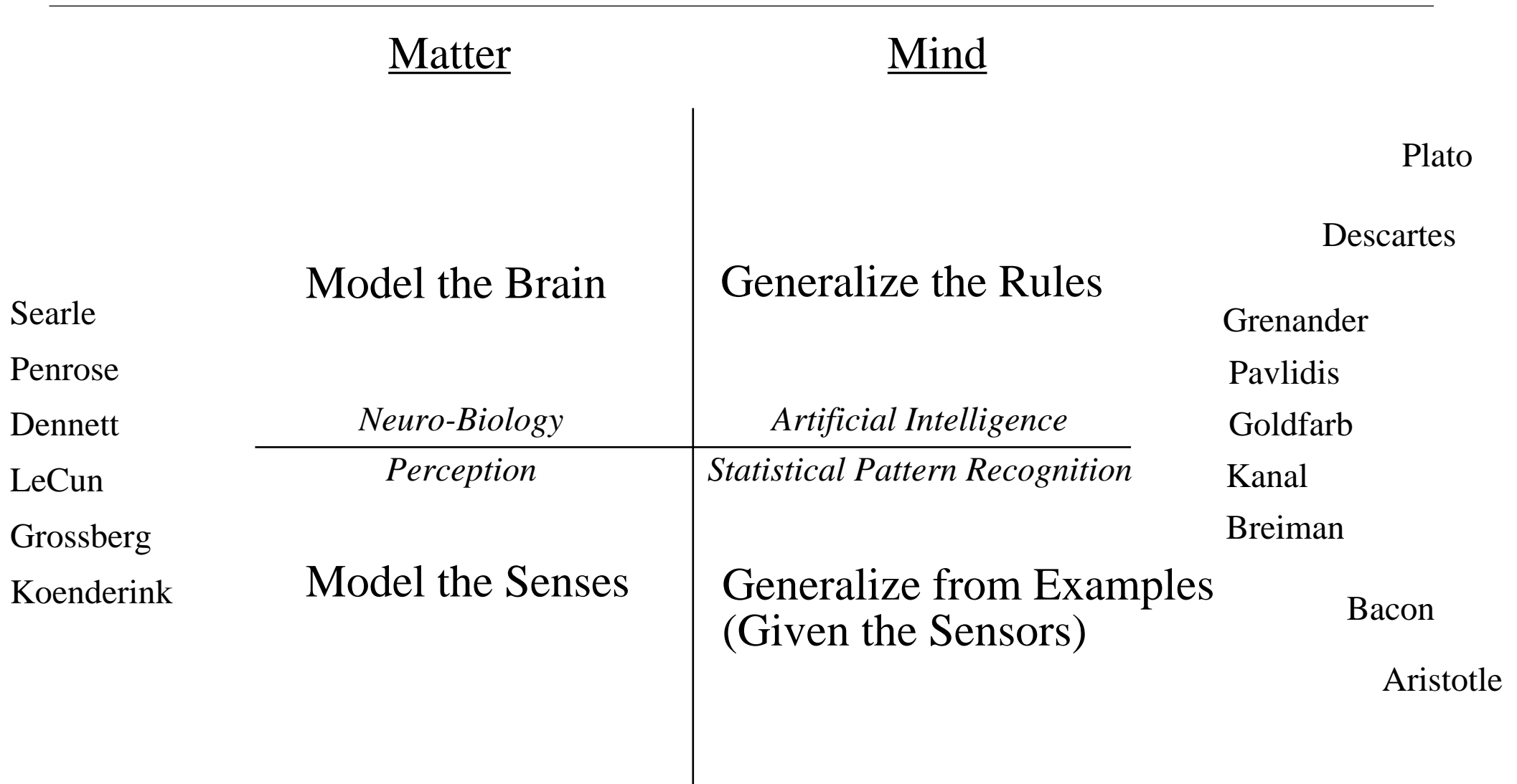
Find the class of a
new, real world object
given a number of examples

Hard AI: The Brain is a Computer



Dennett: There is no scientific need to take into account the concept of consciousness. Its existence can be denied.

Where to Attack the Pattern Recognition Problem?



Shared Weight Network

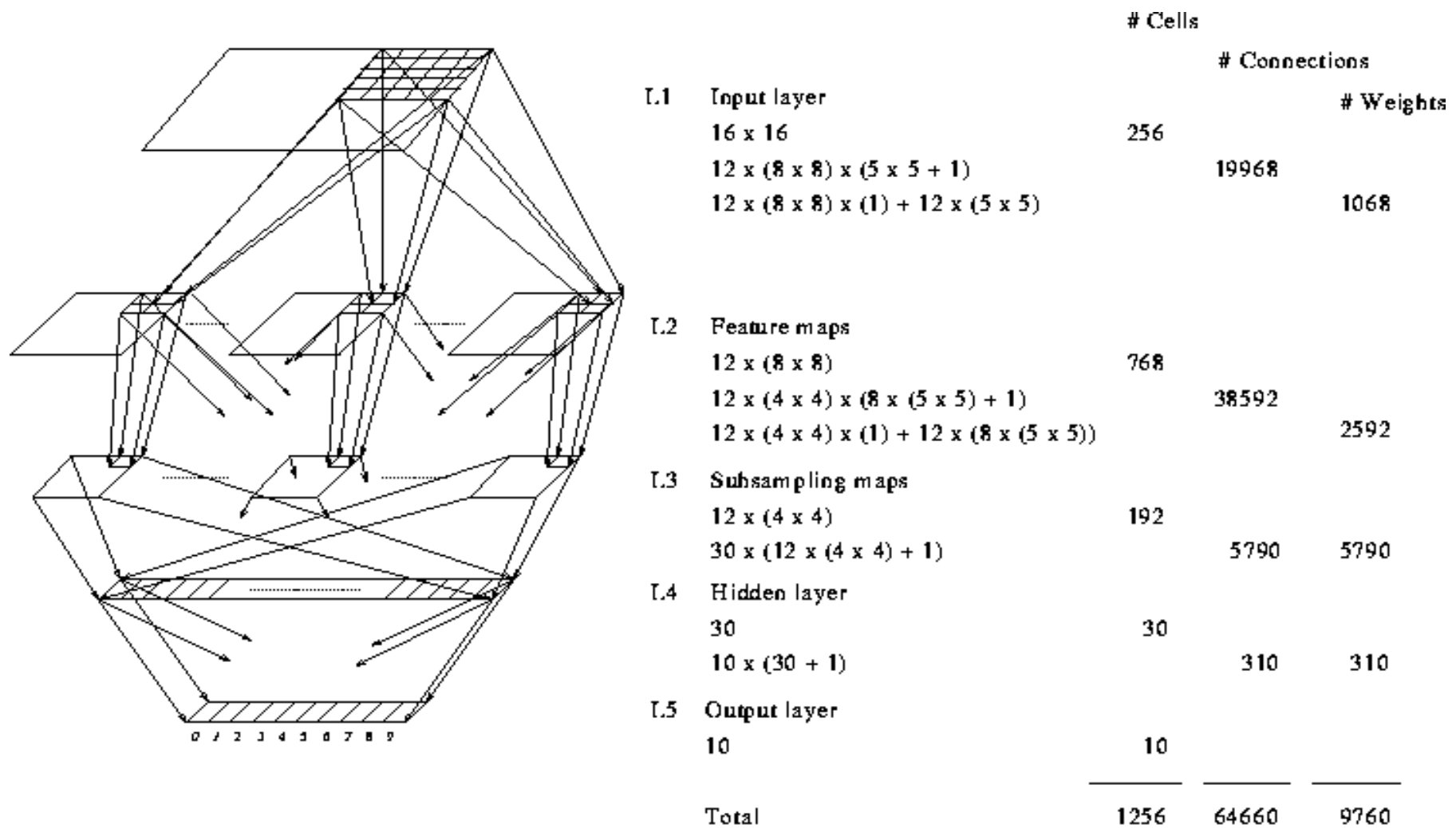
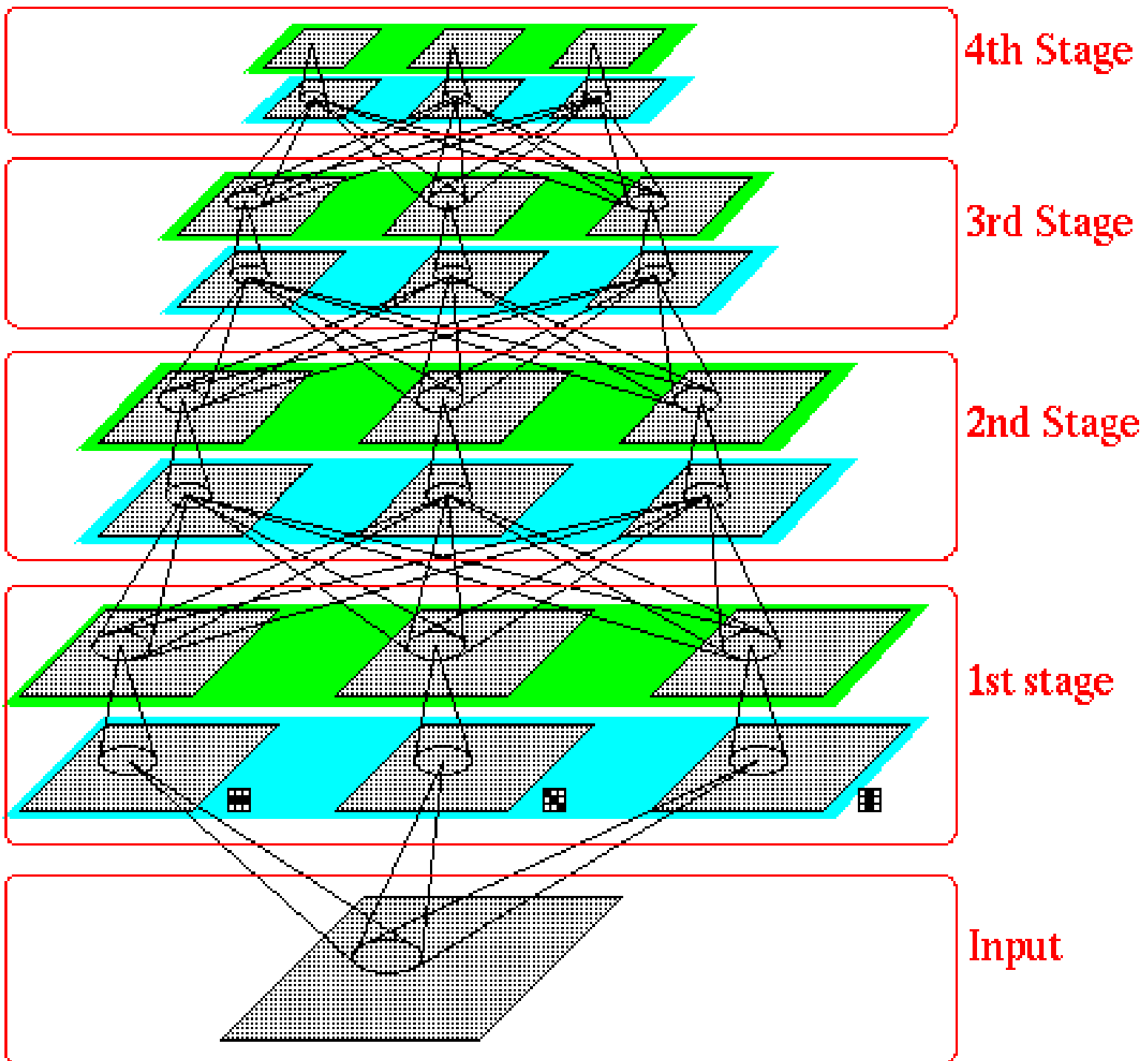


Figure 5.2: An example 2D shared weights ANN.

The Neocognitron Network

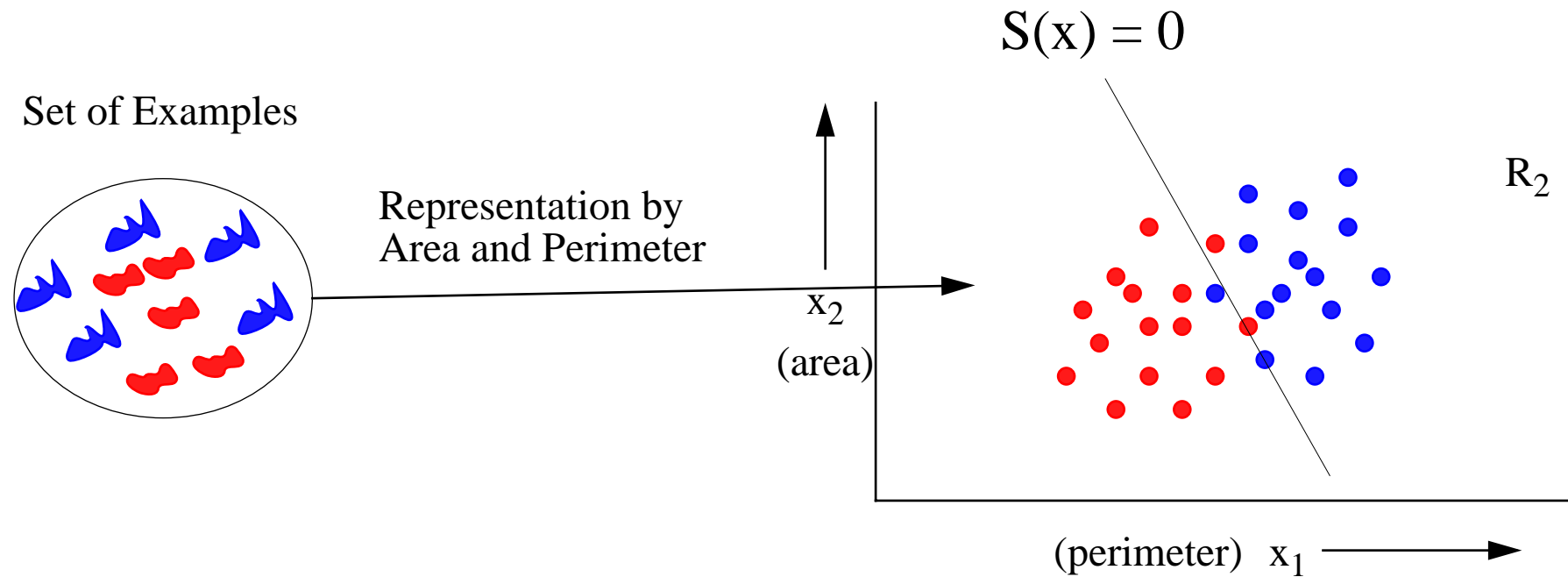


■ — Blurring layer

■ — Feature detecting layer

■ — Cell plane

Sensors, Features and Classifiers



$$\text{Fisher: } S(\mathbf{x}) = (\hat{\mu}_A - \hat{\mu}_B)^T \hat{G}^{-1} \mathbf{x} + \text{constant}$$

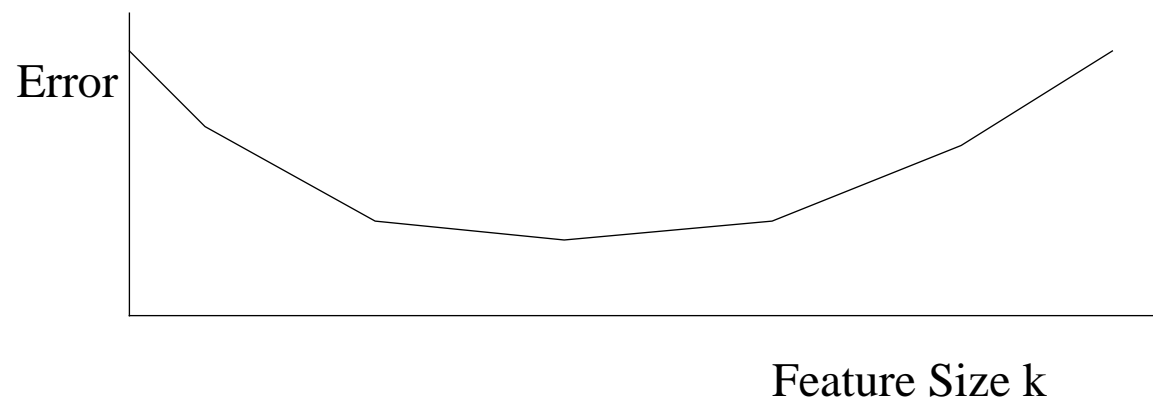
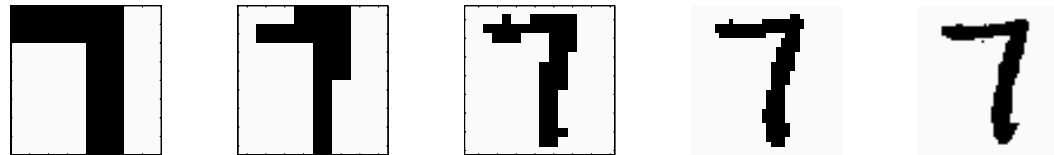
Problems with the Traditional PR Approach

$\mathbf{x} = (x^1, x^2, \dots, x^k)$ - k dimensional feature space

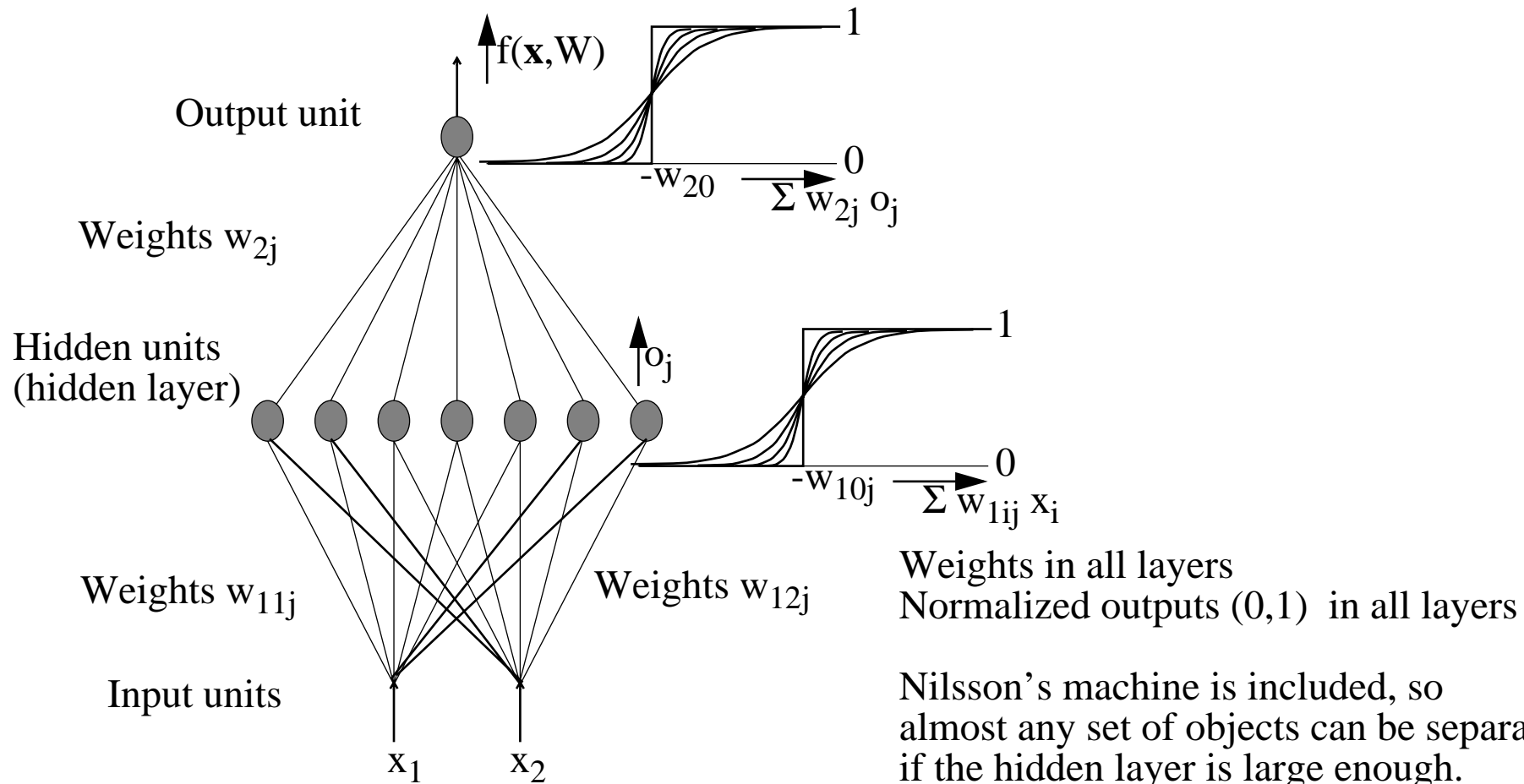
$\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\}$ - training set
 $\{\lambda_1, \lambda_2, \dots, \lambda_m\}$ - class labels

} $D(\mathbf{x})$ - classifier, $\varepsilon = \text{Prob} (D(\mathbf{x}) \neq \lambda(\mathbf{x}))$

$\varepsilon(m)$: monotonically decreasing, $\varepsilon(k)$: peaks !



Neural Networks



More output units are possible
More hidden layers are possible.

Backpropagation Training Rule

Network: $f(\mathbf{x}_p, \mathbf{W})$

Training set: $\{(\mathbf{x}_1, t_1), (\mathbf{x}_2, t_2), \dots, (\mathbf{x}_n, t_n)\}$

Target values (labels): t_p

Network error: $E = \sum_p \{ t_p - f(\mathbf{x}_p, \mathbf{W}) \}^2$

Gradient descent: $\mathbf{W} \leftarrow \mathbf{W} + \Delta\mathbf{W} = \mathbf{W} + \sum_p \Delta_p \mathbf{W}$

Generalized delta rule: $\Delta_p \mathbf{W}$: $\Delta_p w_{ji} = \eta \delta_{pj} o_{pi}, \forall i, j$

Fixed stepsize: η

output units (layer k): $\delta_{pk} = (t_{pk} - o_{pk}) o_{pk} (1 - o_{pk})$

hidden units (layer j): $\delta_{pj} = o_{pj} (1 - o_{pj}) \sum_k \delta_{pk} w_{kj}$

The 'errors' in the lower layers (j) are computed using the corrections of the layers above (k): backpropagation.

The Art of Training a Network

Using a neural network classifier is not straightforward all:

- Architecture (numbers of hidden layers and hidden units)
- input representation
- output representation
- target values
- reject values
- initialization procedure
- batches, partitioned learning set (size?) or individual training
- adding noise (amount?) to input or weights?
- step size η
- momentum term α

"The backpropagation training procedure can be a user's nightmare"

(Weiss & Kulikowski)

Large Network Examples

application	#weights	#samples	error	ref.
text -> speech	25000	5000	0.20	Sejnowski
sonar target rec	1105	192	0.15	Gorman
car control	>36000	1200	car drives on winding road	Pomerleau
back-gammon	>11000	3000	computer champion	Tesauro
sex rec from faces	>36000	90	0.09	Golomb
char rec	9900	5000	0.055	Sato
remote sensing	1800	50	0.05-0.10	Kamata
signature verif.	480	280	0.05	Sabourin

T.J. Sejnowski and C.R. Rosenberg, *NETtalk: a parallel network that learns to read aloud*, The John Hopkins University Electrical Eng. and Comp. Science, 1986.

P. Gorman and T.J. Sejnowski, *Learned Classification of Sonar Targets Using Massively Parallel Network*, IEEE Transactions on ASSP, vol. 36, no. 7, July 1988.

D. Pomerleau, *ALVINN: Ann Autonomous Land Vehicle in a Neural Network*, in: David S. Touretzky, *Advances in Neural Information Processing Systems I*, 1989

G. Tesauro, *Neurogammon wins computer olympiad*, *Neural Computation*, vol. 1, pp 312-323, 1990

B.A. Golomb, D.T. Lawrence, T.J. Sejnowski, *Sexnet: A neural network identifies sex from human faces*, *Adv. in Neural Inf. Proc. Sys. I*, 1989

A. Sato, K. Yamada, J. Tsukumo, and T. Temma, *Neural network models for incremental learning*, *ICNN*, Helsinki, 1991.

S.-I. Kamata, R.O. Eason, A. Perez, and E. Kawaguchi, *A Neural Network Classifier for LANDSAT Image Data*, *Proc. 11th ICPR*, The Hague, Vol 2, 573-576, 1992

R. Sabourin and J-P. Drouhard, *Off-Line Signature Verification Using Directional PDF and Neural Networks*, *Proc. 11th ICPR*, The Hague, Vol 2, 321-325, 1992

Neural Network Appreciation in PR



"Artificial Intelligence and Neural Networks have deceived and spoiled two generations of computer scientists just by these names"
(Rosenfeld, Oulu 1989)



"Neural Networks has brought new enthusiasm and spirit to the next generation of young researchers."
(Kanal, Jerusalem 1994)

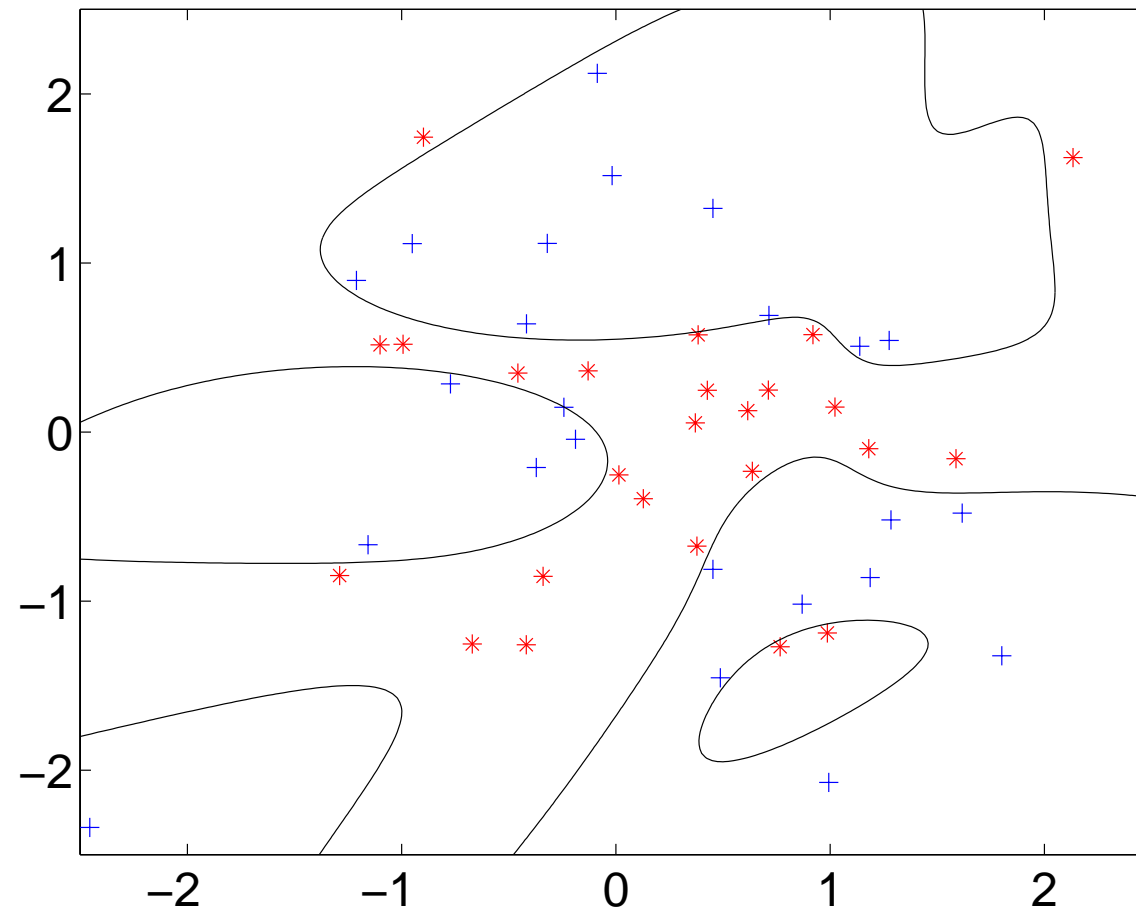
"Just a short look at the architecture of a Neural Network is sufficient to see that the thing simply doesn't have the moral right to show any reasonable performance"
(Breiman, Edinburgh, 1995)

My Neural Network Problem

"Your problem, Dr. Duin, is that you want to understand the neural network. You will have to accept that the interesting aspect of neural networks is that their behavior cannot be understood. " (NN, Delft, 1991)

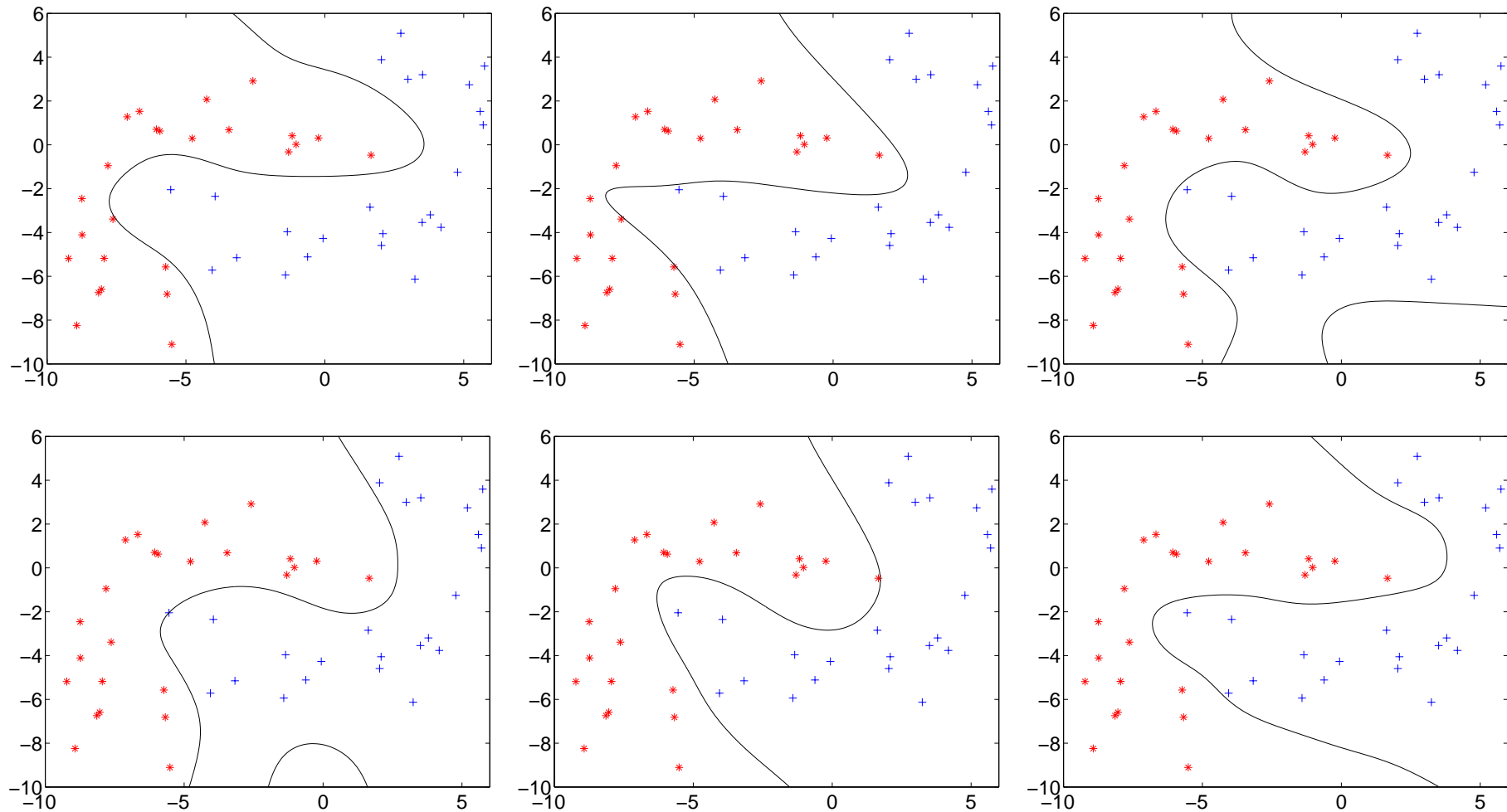
So, it is just magic!

Neural Network as Universal Approximator



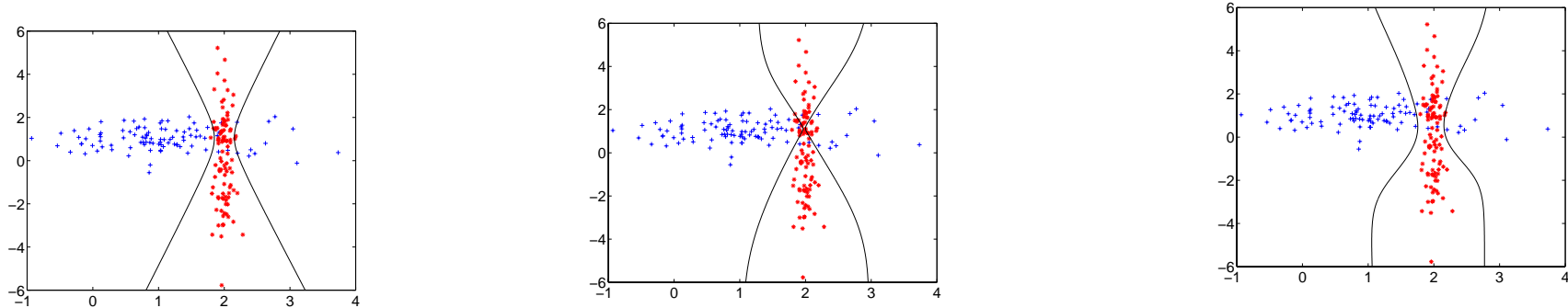
A sufficiently large neural network can solve almost any problem.

Neural Network as Universal Problem Solver



A neural network can solve a problem in many different ways

Training versus Implementing



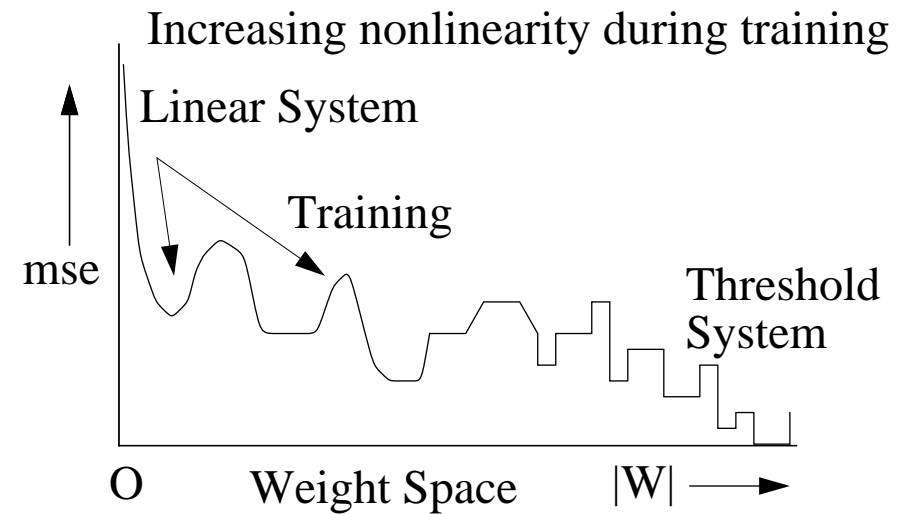
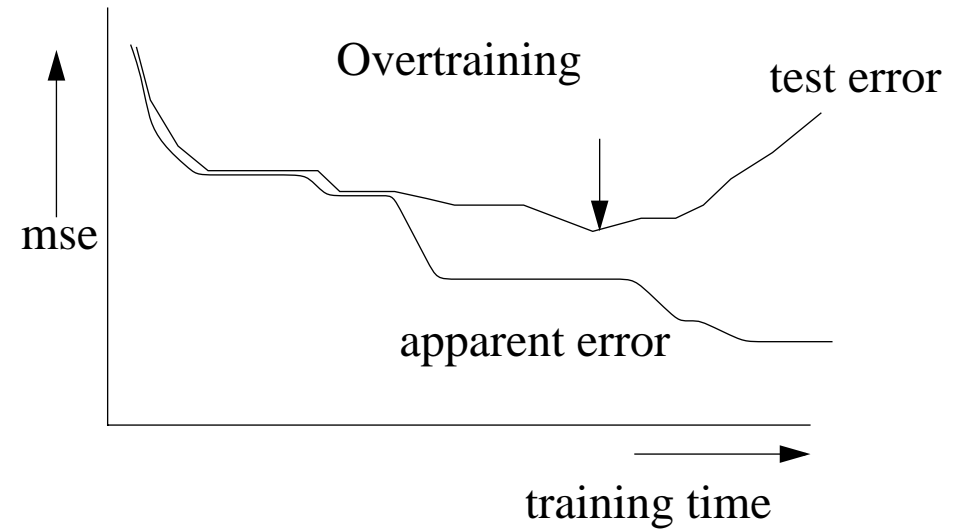
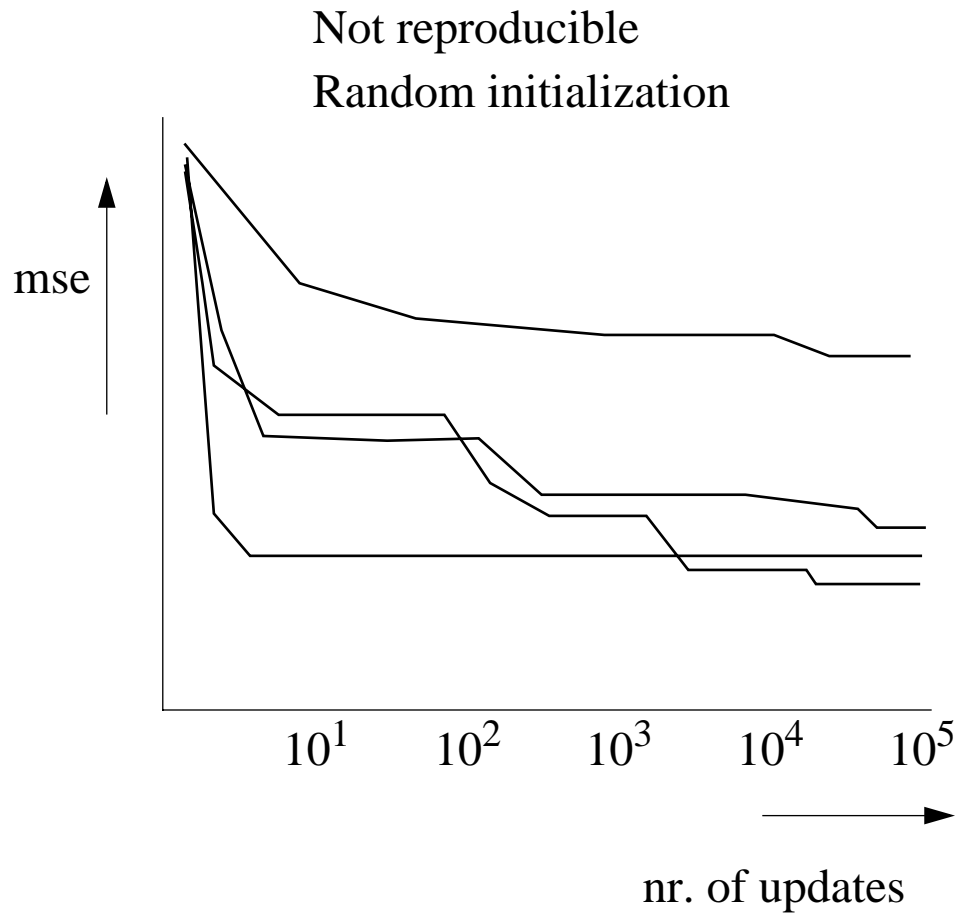
Any function can be implemented on a neural network

Consequently, it can be trained by any rule.

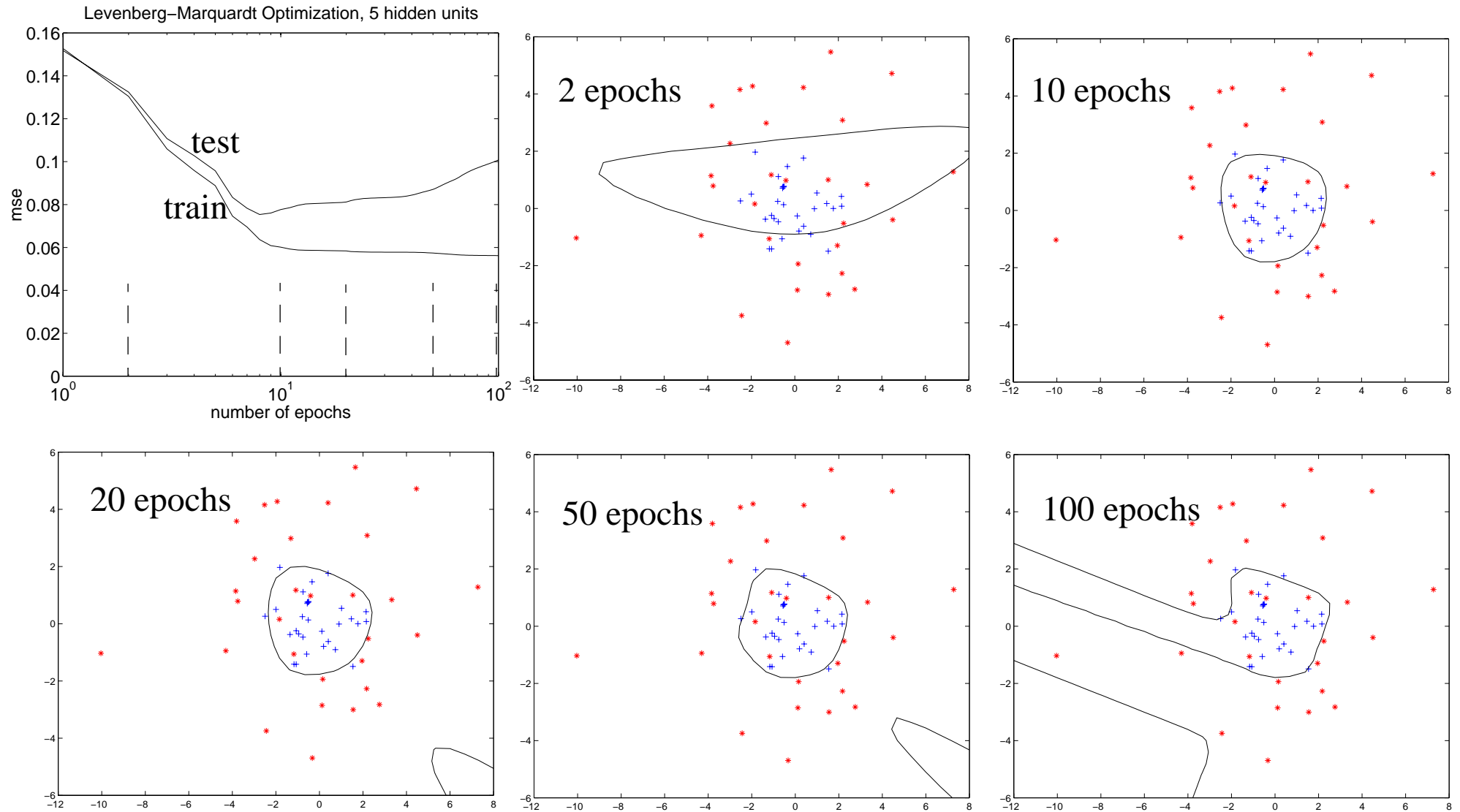
The architecture is general, and thereby not special.

What is special, is the original training rule.

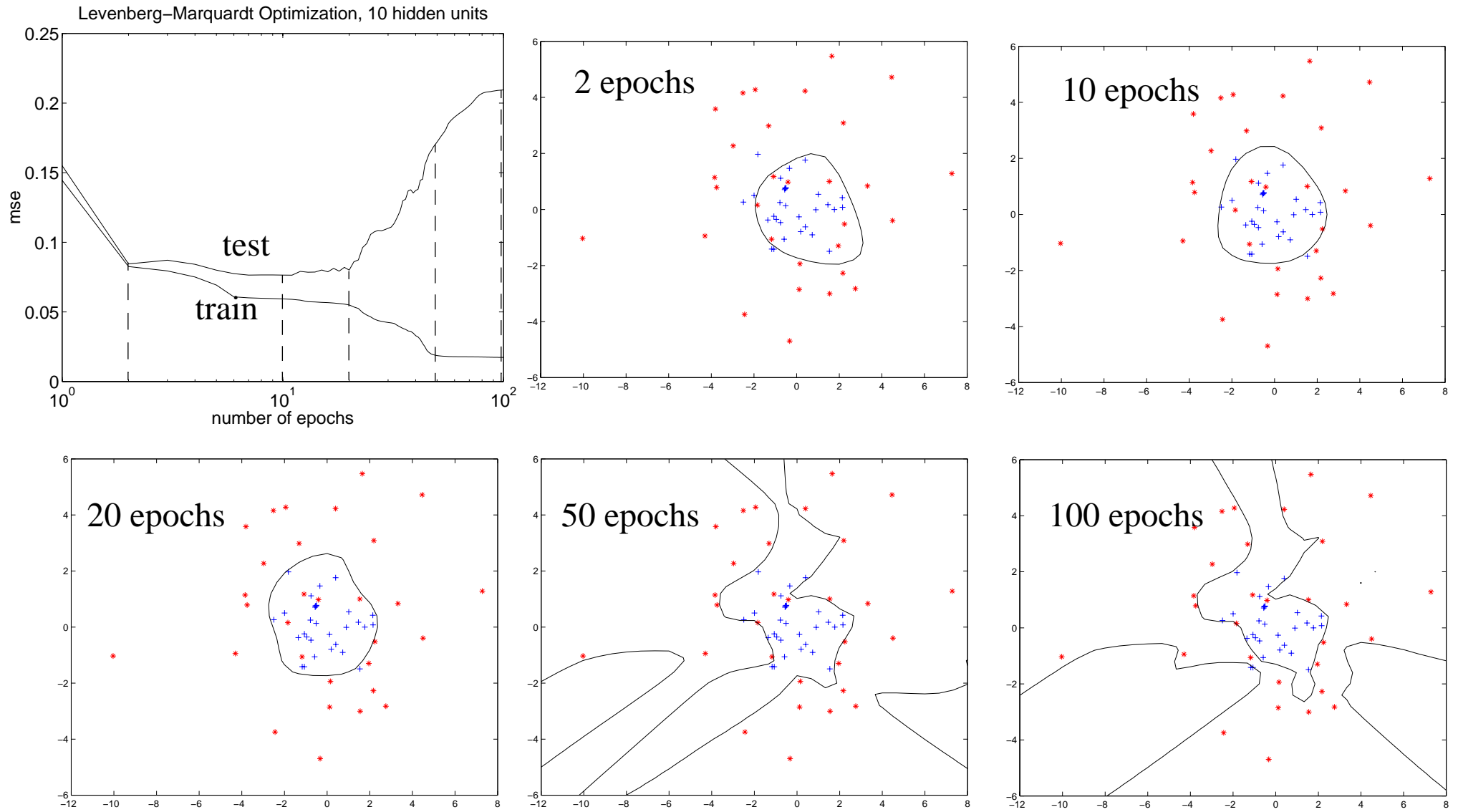
Gradient Descent Training Characteristics



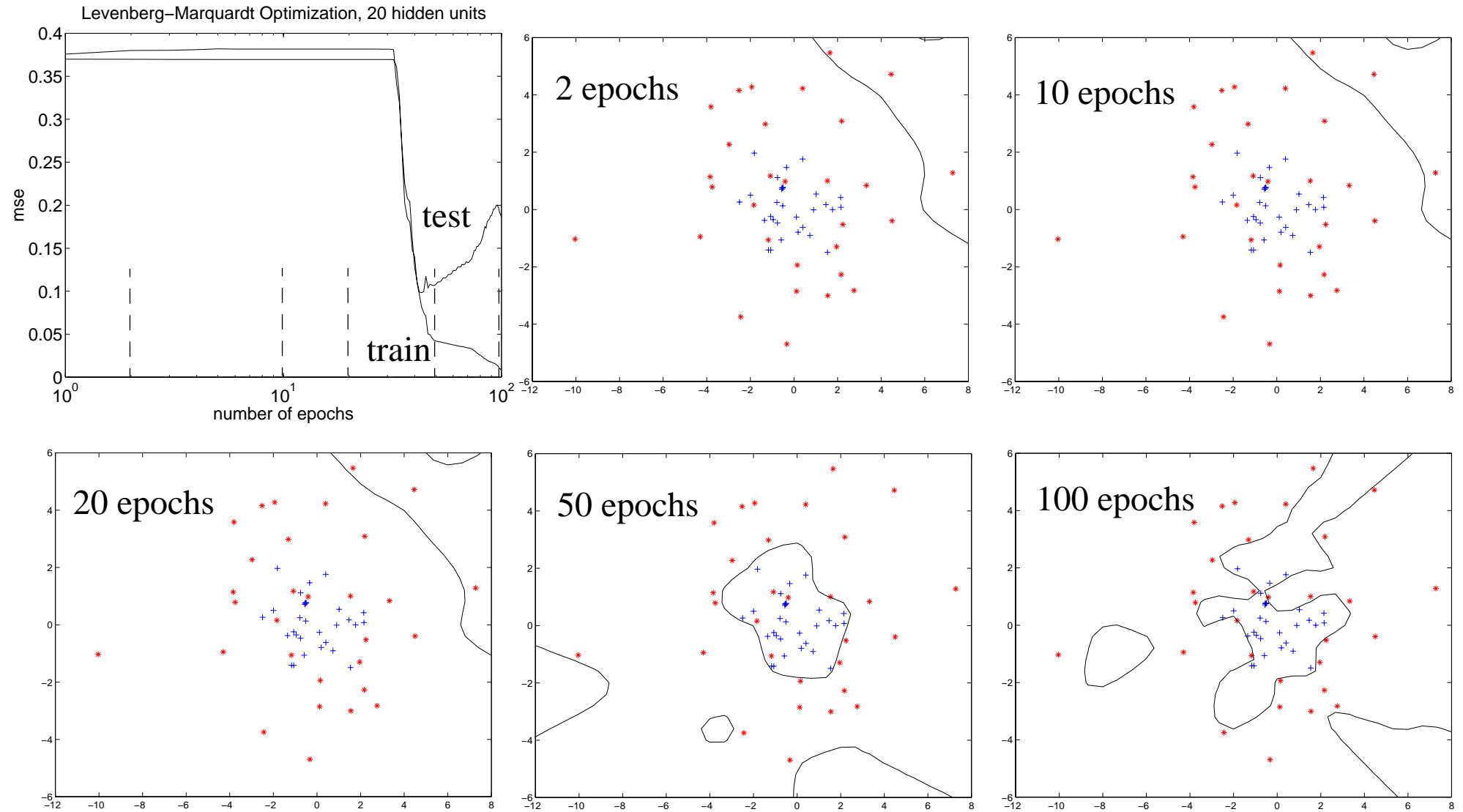
Overtraining Example - 5 Hidden Units



Overtraining Example - 10 Hidden Units



Overtraining Example - 20 Hidden Units



Redundancy

Neural Networks usually have more layers and neurons than necessary.

Training a more simple network, that is able to implement the same function, however, appear to be difficult.

Many neurons may be given random weights (and not trained) without causing problems.

--> Redundancy helps, but demands more time.

Speed and Memory

Training and testing may be very slow and memory demanding.

Special hardware helps, as it does for other procedures.

Are Neural Networks Better?

Neural networks usually are not better
in simple problems

They don't offer fixed procedures.

They offer a complicated toolbox and
not a single off-the-shelf tool.

Application demands a skilled analyst.

Averaged error rates and standard deviations over 10 runs

dataset	NMea n	norm	k-NNR	1-NNR	DTree	ANN
IRIS	<i>0.077</i> 0.019	<u>0.025</u> 0.010	<i>0.048</i> 0.019	<i>0.053</i> 0.017	<i>0.071</i> 0.031	<i>0.052</i> 0.026
IMOX	<i>0.115</i> 0.027	<i>0.102</i> 0.026	0.086 0.018	<u>0.071</u> 0.023	0.092 0.045	0.088 0.031
80X	0.114 0.054	0.123 0.074	<u>0.077</u> 0.083	0.082 0.088	<i>0.255</i> 0.099	0.118 0.078
BLOOD	<i>0.163</i> 0.034	0.125 0.035	0.131 0.035	<i>0.153</i> 0.041	<i>0.158</i> 0.048	<u>0.123</u> 0.033
GLASS	<i>0.569</i> 0.049	<i>0.431</i> 0.098	0.303 0.040	<u>0.286</u> 0.045	<i>0.334</i> 0.052	<i>0.380</i> 0.075
SONAR	<i>0.352</i> 0.072	<i>0.315</i> <i>0.061</i>	0.194 0.050	<u>0.188</u> 0.044	<i>0.307</i> 0.043	<i>0.236</i> 0.034
DNORM	<i>0.334</i> 0.045	<i>0.151</i> 0.041	<i>0.185</i> 0.045	<i>0.212</i> 0.038	<i>0.344</i> 0.045	<u>0.121</u> 0.017

Conclusions on Neural Networks

Neural Networks are oversized, general function approximators that work because of the training rule:

- start from linearity.
- stop in one of the first moderately nonlinear local minima.
- have many (built-in) regularization possibilities, including a slow optimization rule (back-propagation).
- do not reproduce.
- are computational intensive.

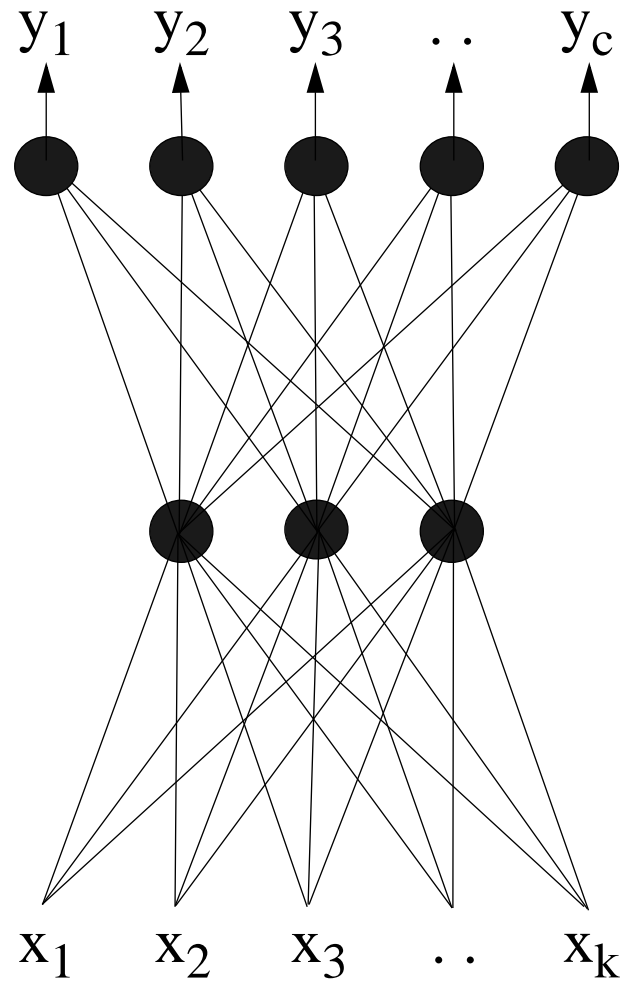
Final Conclusion on Neural Networks for Pattern Recognition

<u>Matter</u>	<u>Mind</u>
Model the Brain	Generalize the Rules
<i>Neuro-Biology</i>	<i>Artificial Intelligence</i>
<i>Perception</i>	<i>Statistical Pattern Recognition</i>
Model the Senses	Generalize from Examples (Given the Sensors)

Pattern Recognition is looking for general procedures for learning from examples.

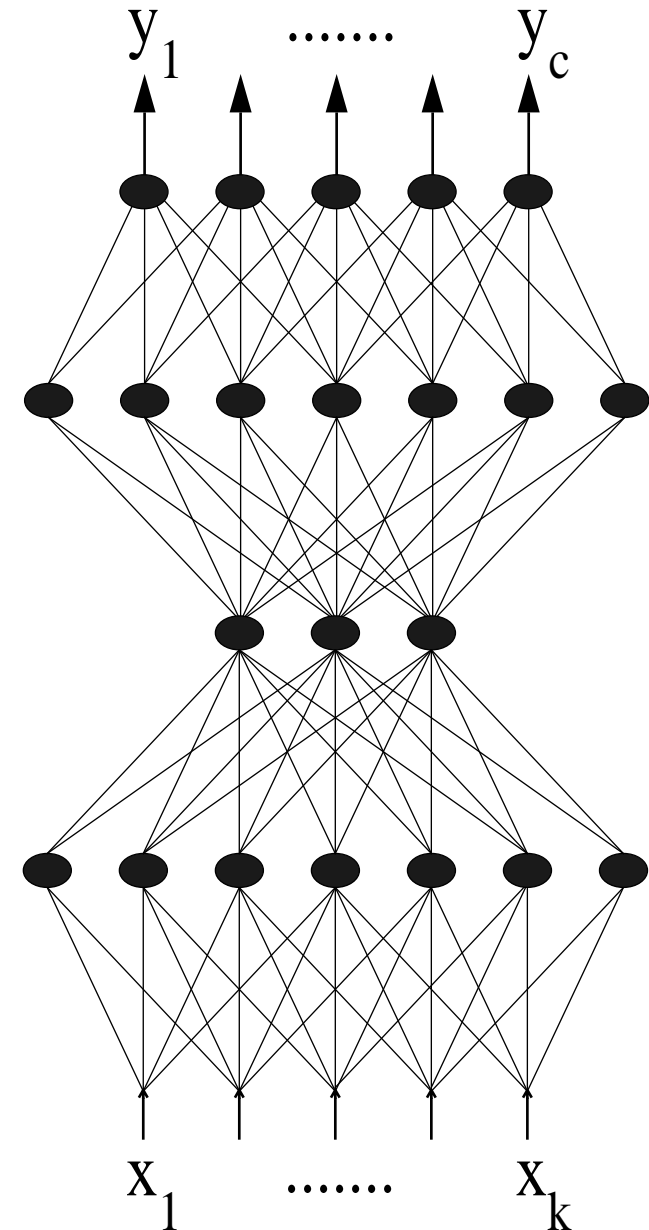
Neural Networks offer a toolbox to find specific solutions for specific problems.

Nonlinear Mapping

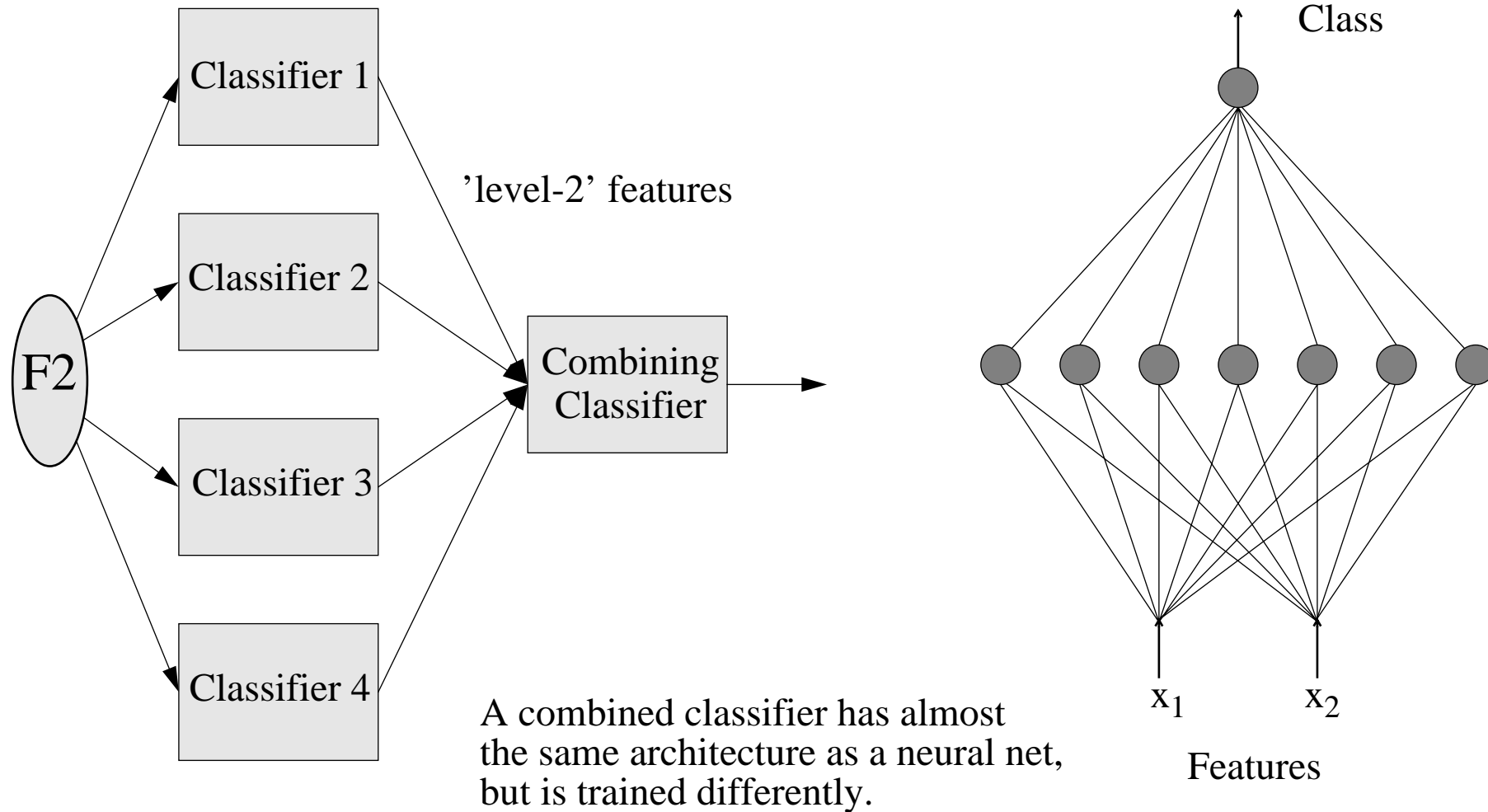


Interpretation:

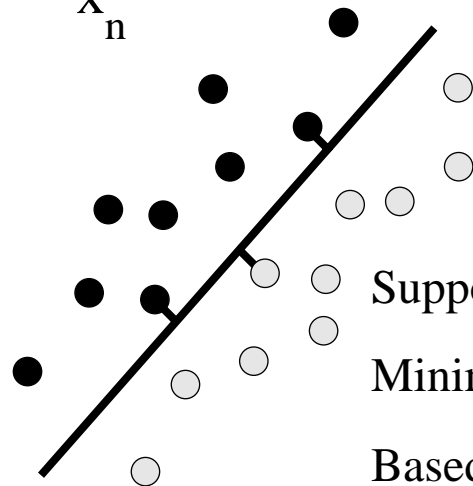
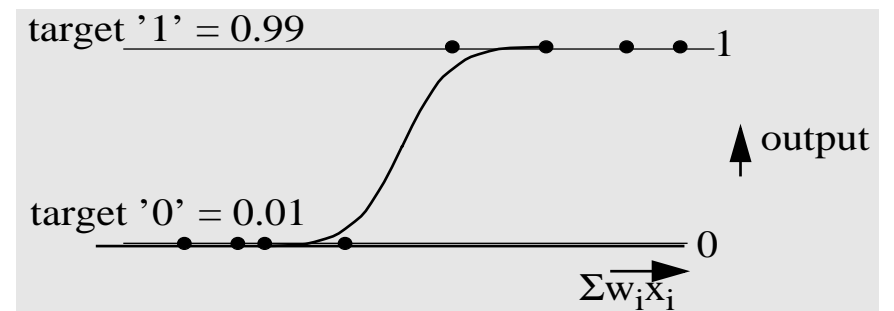
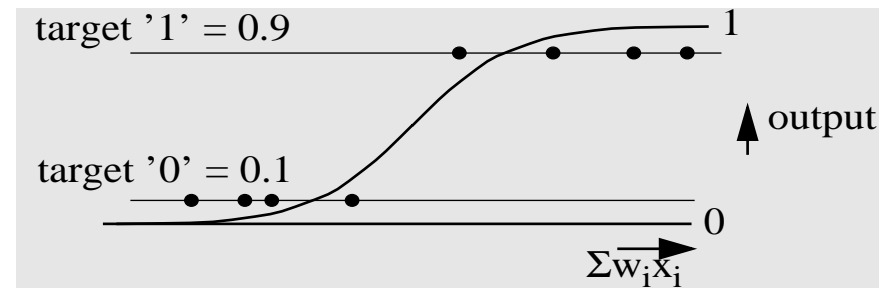
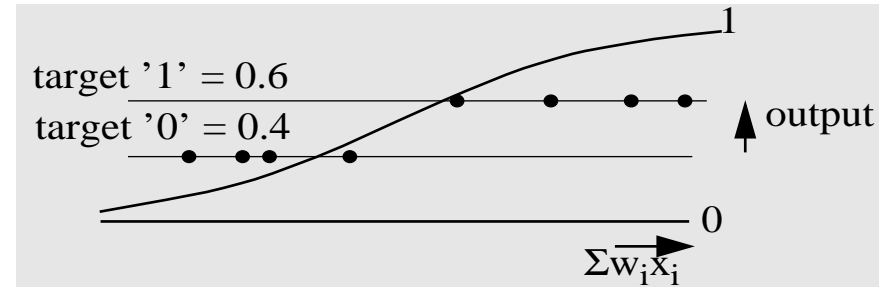
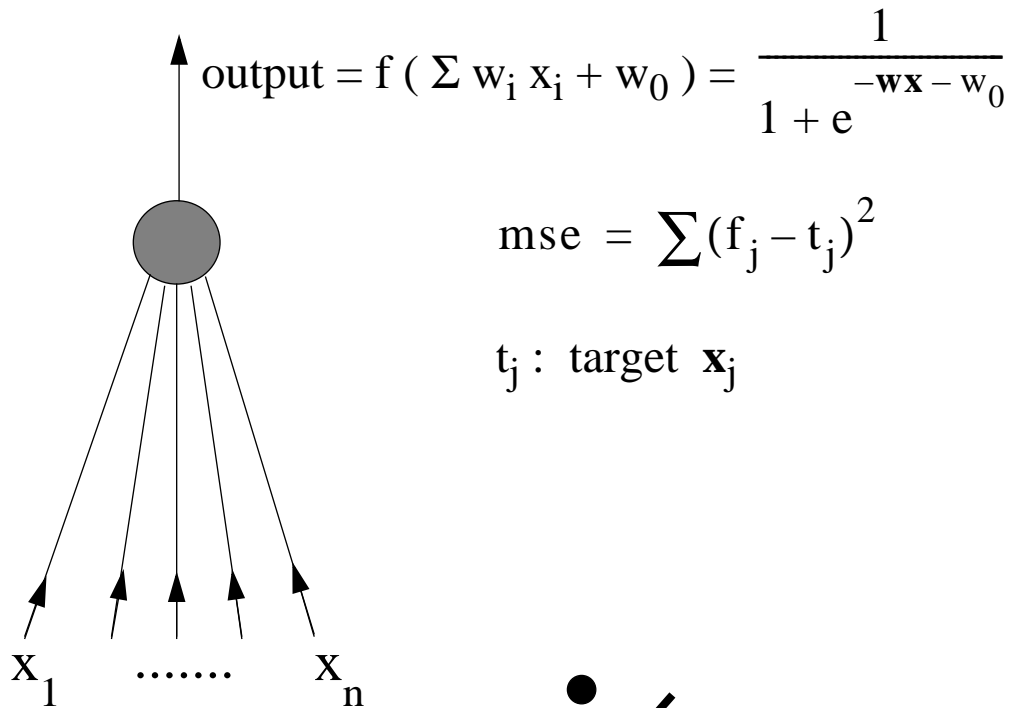
1. Find a (non)linear subspace.
2. Classify.



Combining Classifiers



Neural Network --> Support Vector Machines



Support Vector Classifier

Minimize training set to a support set

Based on inner products $X_i^T X_j$

What Has Been Learned

Understanding the **use of redundancy**, oversized systems and **regularization**.

The use of controlled moderate **nonlinearities**: Nonlinear mapping techniques.

Better classifiers: support vector machines

Soft outputs: confidences and fuzzy memberships

The construction and use of **complicated systems**: combined classifiers.