# An Introduction to Neural Networks

Reference: B.J.A. Kröse and P.P. van der Smagt (1994): An Introduction to Neural Networks, Poglavja 1-5, 6.1, 6.2, 7-8.

- B.J.A. Kröse and P.P. van der Smagt (1994): An Introduction to Neural Networks, Seventh edition, The University of Amsterdam.
- J. Sjöberg et al. (1995):Non-linear Black-Box Modeling in System Identification: a Unified Overview, Automatica, Vol. 31, 12, 1691-1724.
- M. Agarwal(1997): A systematic classification of neural-networkbased control, IEEE Control Systems Magazine, Vol. 17, No. 2, 75-93.
- M. Nørgaard (1995): Neural Network Based System Identification Toolbox User's Guide, Technical Report 95-E-773, Institute of Automation, Technical University of Denmark, Lyngby.

http://www.iau.dtu.dk/research/control/nnsysid.html

M. Nørgaard (1995):Neural Network Based Control System Design Toolkit User's Guide, Technical Report 96-E-830, Institute of Automation, Technical University of Denmark,Lyngby. http://www.iau.dtu.dk/research/control/nnctrl.html

#### Introduction: historical overview

#### **Milestones**

• 1943 – dawn of Artificial Neural Networks (ANN) McCulloch in Pitts – formal representation of neuron, mathematical abstraction – a model - of biological neuron





 $y_i(t) = f(\sum_i w_{ij}(t) x_j(t) + \theta_i(t))$ 

3

The neuron of McCulloch and Pitts

- Treshold activation function: elements of logic functions
- Conections between units into networks
- Formal model of nevron unchanged until present days

- 1949 Hebb: psychologist, the first learning rule
- **1959 Rosenblatt:** the first who used "*perceptron*" for single layer network
- **1960** Widrow, Hoff: *adaline* (ADAptive LINear Element) single layer, the first *analiticaly* derived learning rule (least squares), until then only heuristical learning, in analogy with BIO systems
- 1963 Widrow, Smith: inverted pendulum, adaline

#### n Perceptron (1959)

- <sup>1</sup> Activation function: switching function
- <sup>1</sup> Single layer network
- <sup>1</sup> Can not represent XOR function
- n Adaline (Adaptive linear element 1960)
  - <sup>1</sup> Activation function: linear function
  - Learning with "delta rule" (least squares optimisation method)
  - 1 linear models

Explanations

n Regressionn Classification

n Learning: optimising weights of ANN

Systems modelling from data

# Learning methods

#### n Hebbs rule

$$\Delta w_{ij} = \gamma y_i x_j$$

n Delta rule (Widrow-Hoff rule)

$$\Delta w_{ij} = \gamma x (d_i - y_i) x_j$$

7

n Least squares optimisation method

- 1949 Hebb: psychologist, the first learning rule
- 1959 Rosenblatt: the first who used "perceptron" for single layer network
- 1960 Widrow, Hoff: *adaline* (ADAptive LINear Element) single layer, the first *analiticaly* derived learning rule (least squares), until then only heuristical learning, in analogy with BIO systems
- 1963 Widrow, Smith: inverted pendulum, adaline
- **1969** Minsky, Papert show the limitations of perceptrona, that can not be used for classification of elements that are not linearly separable -XOR function can not be represented
  - nothing much happens until app. 1982
    The end of single layer networks period!
- 1986 Rumelhart and co.: "backpropagation" learning rule for multi-layer perceptron classification of elements that are not linearly separable, the nonlinear mapping!!!
  The renaissance of ANN

# Multi-layer perceptron



 $y_{i} = f_{i} \left( \sum_{i} w_{ij} f_{j} \left( \sum_{k} w_{jk} x_{k} + w_{0k}^{1} \right) + w_{0j}^{2} \right)$ 





Systems modelling from data

#### Radial basis function (RBF) network

• all weights between inputs and hidden layer  $\equiv 1$ 



# n Back-propagation method (BP)

- n multi-layer feed-forward networks
- n delta rule generalised for nonlinear problems (gradient optimisation method)
- n BP algorithm is used for the calculation of cost function gradienta
- n Applicable for all sorts of nonlinear systems, but differences in computation complexity
- n improvements:
  - n Learning rate with momentum
  - n Learning per pattern

ANN topology
 Feed-forward networks

S Recurrent networks

#### **S** ANN learning

- Supervised (associative) learning
- S Unsupervised learning (self-organisation)

# S ANN application

- S Regression
- S Classification

# Recurrent networks

- n The main difference from feed-forward networks is that they contain feedback connections
- n Mainly self-organised networks
- n Rekursive neural networks
- n Hopfield neural network
- n Boltzmann machines

#### Hopfield neural network

- Signum activation function or linear function with saturation
- n Feed-back connections
- n Hopfield network as associative memory
- <sup>n</sup> Theoreticaly intriguing, but less interesting for practice

### **Boltzmann machines**

- n Hopfield network with hidden layers
- Stochastic update rule instead of deterministic one principle of annealing

# Self-organised neural networks

- Instead of input/output data pairs as in supervised learning, self-organised NN use only input data
- n Tipically used for classification: clustering, vector quantification, dimensionality reduction, feature extraction (pattern recognition)
   n Kohonen network

- 1949 Hebb: psychologist, the first learning rule
- 1959 Rosenblatt: the first who used "perceptron" for single layer network
- 1960 Widrow, Hoff: adaline (ADAptive LINear Element) single layer,
- the first analiticaly derived learning rule (least squares),
- until then only heuristical learning, in analogy with BIO systems
- 1963 Widrow, Smith: inverted pendulum, adaline
- 1969 Minsky, Papert show the limitations of perceptrona, that can not be used
- for classification of elements that are not linearly separable -XOR function
- can not be represented
- - nothing much happens until app. 1982
- The end of single layer networks period!
- 1986 Rumelhart and co.: "backpropagation" learning rule for multi-layer
- perceptron classification of elements that are not linearly separable,
- the nonlinear mapping!!!
- 1988 Psaltis: ANN used as controller
- 1990 Narendra in Parthasarathy: system identification and control with ANN
- 1995 Sjöberg et co.: nonlinear system identification nonlinear regression

n Newton optimisation methods (2nd order gradient optimisation method) n Gauss-Newton modification n Levenberg-Marguardt modification n Other kinds of learning (various determenistic and stochastic optimisation methods)





• **Definition of ANN:** does not exist, only *common main properties*:

- composed of large number of simple and interconnected elements
- in general adaptive topology

interconnections change – learning rules change "weights" number of network elements is changing

- topology enables parallel information processing (?)
- data are processed according to ANN states and inputs
- Have ANN fulfilled expectations?
- *Early days*: comprehension of BIO systems with mathematical modelling,

strong connection between BIO research and ANN, large expectations



#### **Applications of ANN**

- n Aerospace (automatic pilots, fault detection etc.)
- n Automotive industry (automatic control etc.)
- n Finance sector (document recognition etc.)
- n Electronics (control, pattern recognition etc.)
- n Artificial speach (source recognition etc.)
- n Production systems (control, forecasting etc.)
- n Medicine (signal processing applications etc.)
- n Accountancy (trends forecasting etc.)
- n Robotics (control, signals processing etc.)
- n Telecommunications (data compression etc.)
- n Transportation sciences (diagnosis systems etc.)
- n Military industry (radar and sonar signals processing etc.)
- n Entertainment industry (animation, special effects etc.)
- n Insurance (assets optimisation etc.)

- ANN – **intelligent systems** ??? pattern recognition systems learned from examples generalisation ability

*Paper 1992* - (Kohonen, Neural Networks, 1988): "ANNs are parallel inter-connected networks of simple computational elements <u>which are intended</u> to interact with the objects of the real world in a similar way to biological nervous <u>systems</u>."

ANN development in two separate directions:

 classification, pattern recognition – very close to ANN
 regression (ANN for systems identification and control, time-series identification )

 differences: methodology, topology, learning rules, fields of applications

#### **Important properties of ANN**

a) ANN as nonlinear mapping – at least two network topologies that are *universal approximators* unknown function  $f(\underline{x}), f(\underline{x}): \mathfrak{R}^{nx} \to \mathfrak{R}^{ny}, \underline{x} \circ \underline{y}$ *ANN* approximates  $f(\underline{x})$ 

b) Ability to learn from examples  $(\underline{v}, \underline{y})$ ANN parameter optimisation based on selected cost function, e.g.

$$E(k) = \frac{1}{2} \sum_{j=1}^{n} (d_j(k) - y_j(k))^2$$

with selected optimisation algorithm (learning rule)

c) Generalisation ability (comes from a) )

# Multi-layer perceptron





# Sigmoid or tanh activation function

$$y = g(\mathbf{u}, \mathbf{w}) = \sum_{j=1}^{N_h} c_j a \left( \sum_{i=1}^{N_i} w_{ij} u_i + b_j \right) + d$$

# Multi-layer perceptron – Example 1



# Multi-layer perceptron – Example 2



# The Radial Basis Function Network





# RBF network – Example

