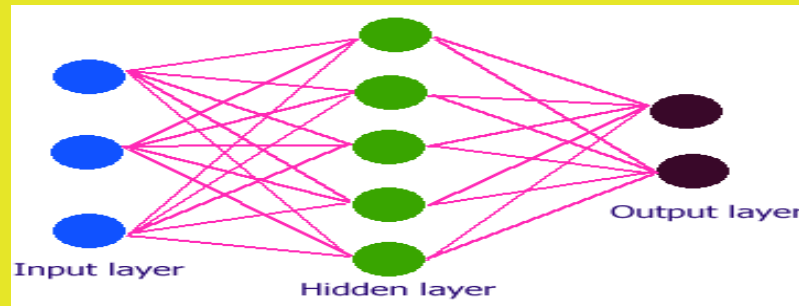


Statistical Data Mining



Artificial Neural Networks

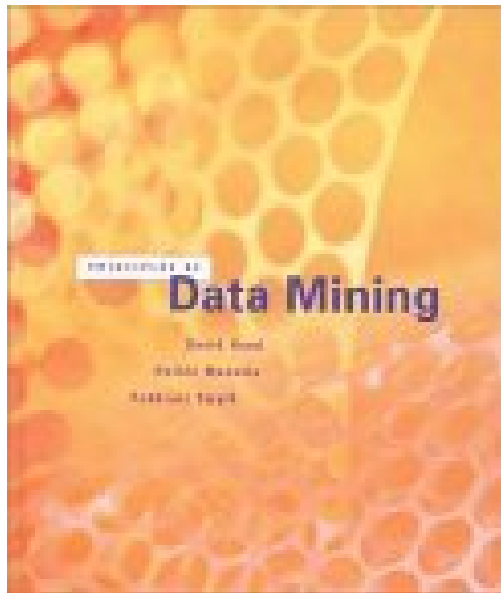
Professor Dr. Gholamreza Nakhaeizadeh

Content

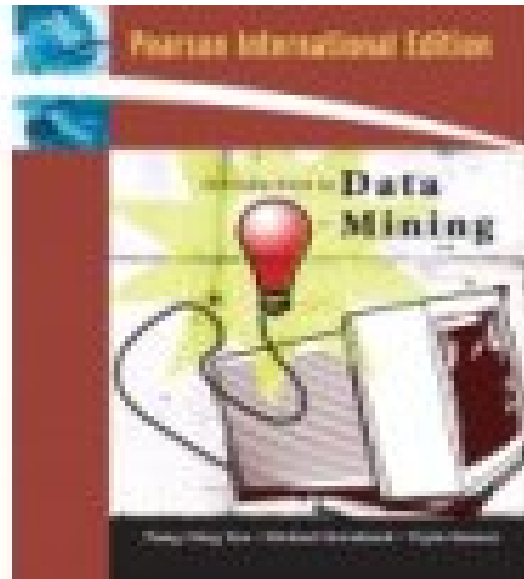
Introduction

- Rise and fall of Neural Networks
- Input function of neuron
- Activation function of neuron
- Output (function) of neuron
- Feed-Forward networks
- Feedback networks
- Learning process
- Coding and decoding methods
- Perceptron
- Backpropagation
- Weakness and Strength of ANN

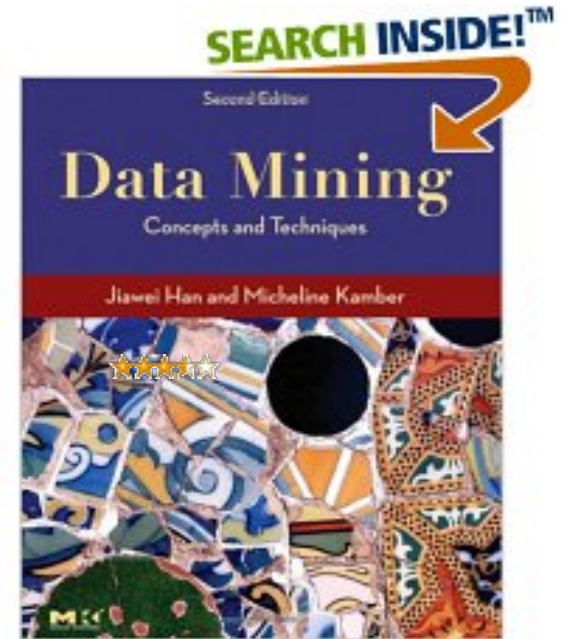
Literatur used (1)



Principles of Data Mining
[David J. Hand](#), [Heikki Mannila](#),
[Padhraic Smyth](#)



Pang-Ning Tan,
Michael Steinbach,
Vipin Kumar



[Jiawei Han](#) and
[Micheline Kamber](#)

Literature Used (2)

<http://cse.stanford.edu/class/sophomore-college/projects-00/neural-networks/>

<http://www.cs.cmu.edu/~awm/tutorials>

<http://www.crisp-dm.org/CRISPwP-0800.pdf>

http://en.wikipedia.org/wiki/Feedforward_neural_network

http://www.doc.ic.ac.uk/~nd/surprise_96/journal/vol4/cs11/report.html#Feedback%20networks

<http://www.dmreview.com/>

<http://www.planet-source-code.com/vb/scripts/ShowCode.asp?lngWId=5&txtCodeId=378>

http://download-uk.oracle.com/docs/html/B13915_02/i_olap_chapter.htm#BABCBDFA

http://download-uk.oracle.com/docs/html/B13915_02/i_rel_chapter.htm#BABGFCFG

<http://training.inet.com/OLAP/home.htm>

<http://www.doc.gold.ac.uk/~mas01ds/cis338/index.html>

<http://www.maths.anu.edu.au/~steve/pdcn.pdf>

www.kdnuggets.com

~~The Data Warehouse Toolkit~~ by Ralph Kimball (John Wiley and Sons, 1996)

Building the Data Warehouse by William Inmon (John Wiley and Sons, 1996)

Studenmund, A. H. (2006). Using Econometrics, A practical Guide. Pearson International Edition

Artificial Neural Networks

Data Mining Algorithms

Machine Learning

- Rule Based Induction
- Decision Trees
- Neural Networks
- Conceptual clustering
-

Statistics

- Discriminant Analysis
- Cluster Analysis
- Regression Analysis
- Logistic Regression Analysis
-

Database Technology

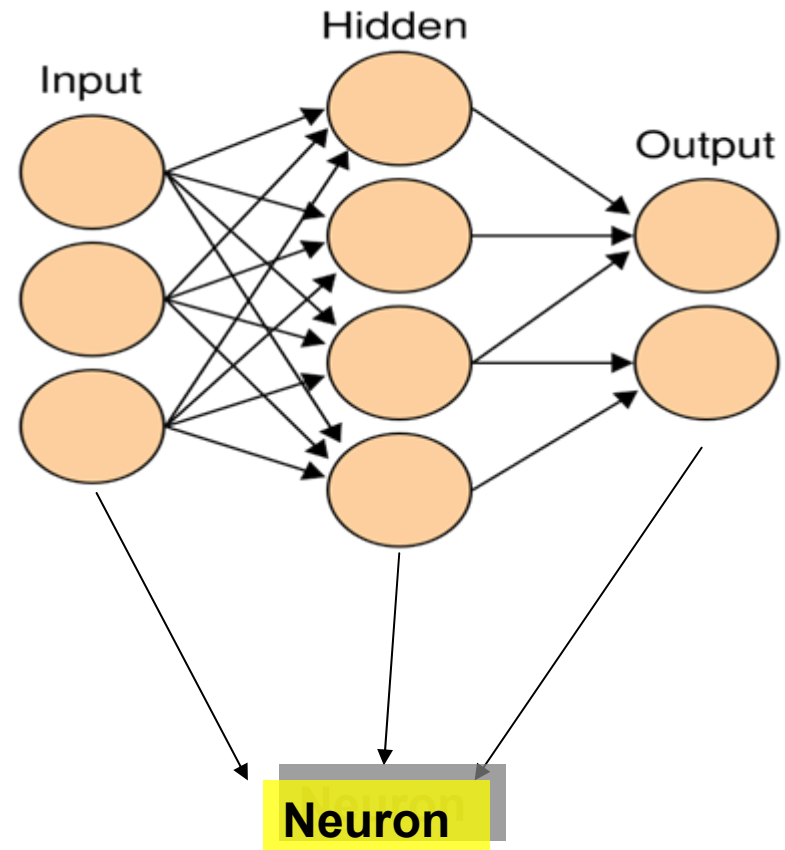
- Association Rules
-

Artificial Neural Networks

Artificial Neural Networks (ANN)

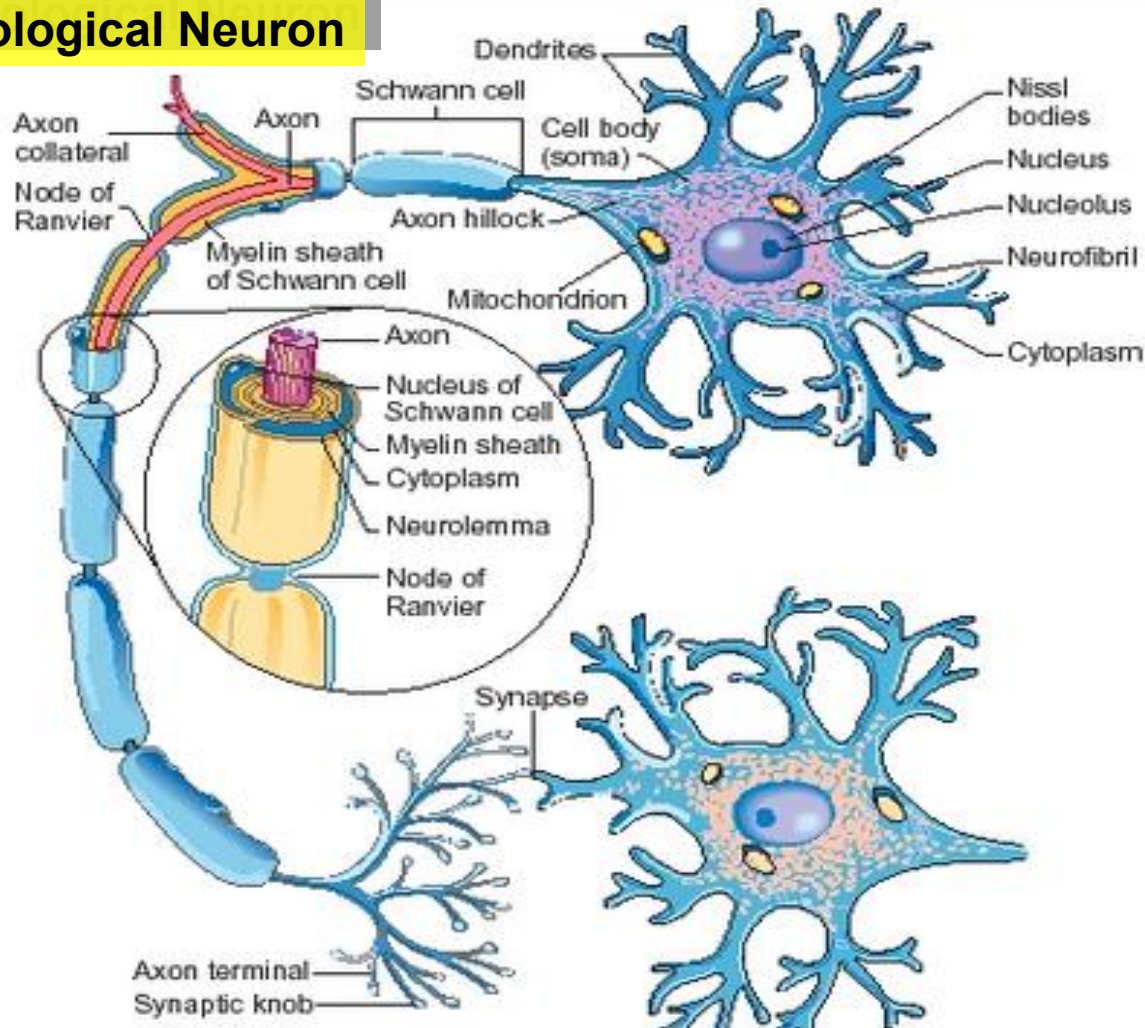
Introduction

- Inspired by the way biological nervous systems e.g. the brain, process information
- An ANN consist of an input layer, an output layer and one or more hidden layer(s):
 - exchanges and process information
 - has learning capability
 - is an adaptive system that changes its structure during the learning phase



Artificial Neural Networks

Biological Neuron

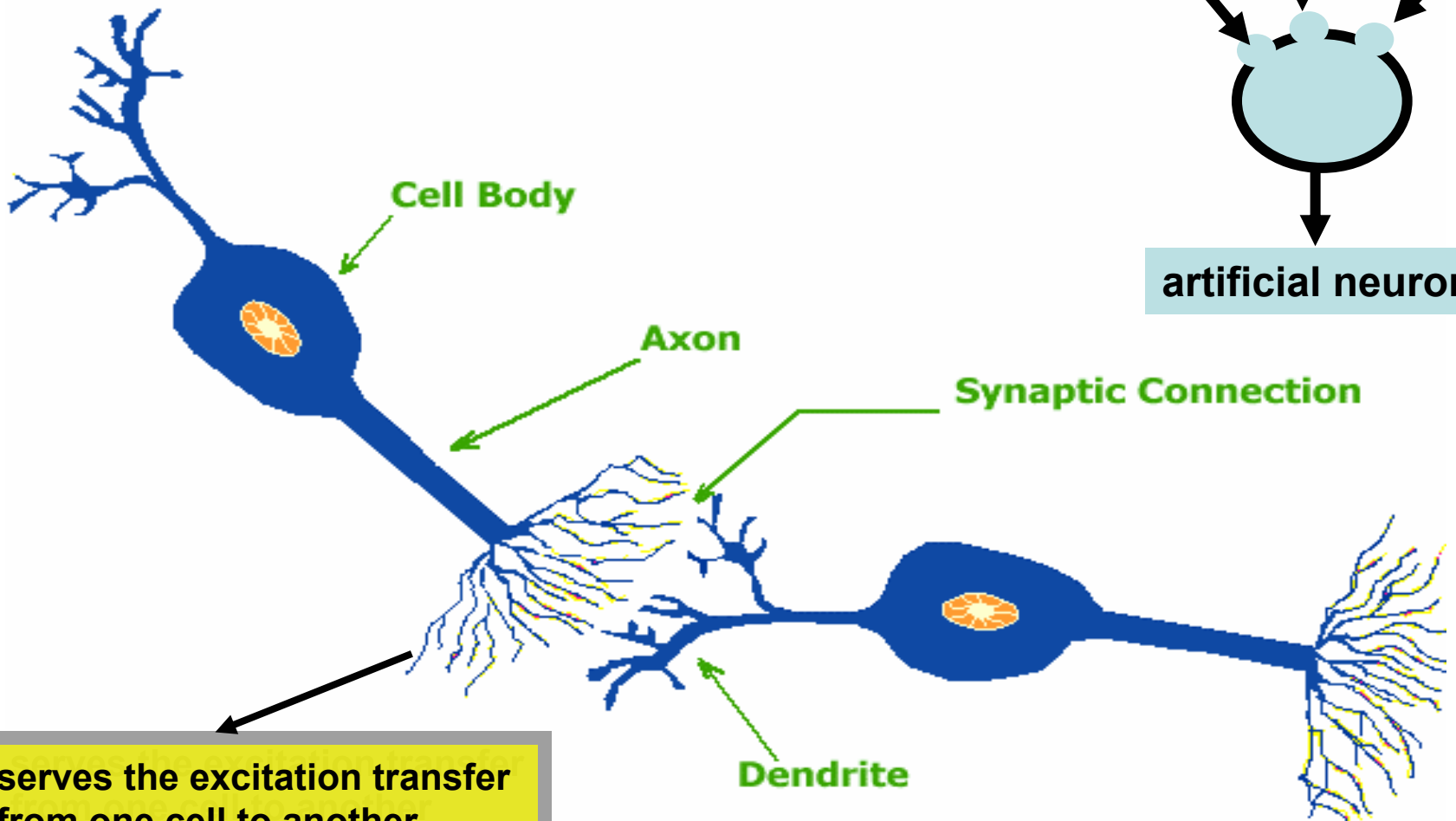


Neuron. (Diagram by Hans & Cassidy. Courtesy of Gale Group.)

Artificial Neural Networks

Biological Neuron

Source: <http://www.intelligentsolutionsinc.com/part1.htm>
great brain



Artificial Neural Networks

Rise and fall of Neural Networks

Fall

- In 1969 Minsky and Papert showed the limitations of single layer Perceptrons
- Their results caused that a lot of researchers loose their interest

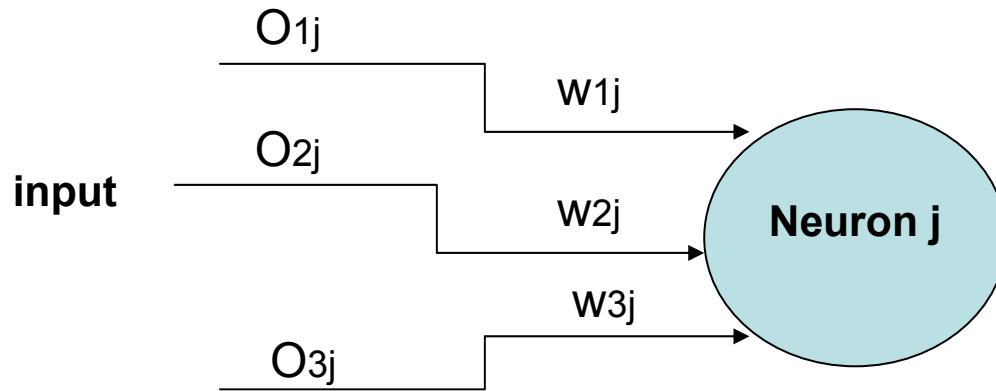
Rise

- In the 70's and 80's, it was shown that multilevel perceptrons don't have These shortcomings
 - Paul J. Werbos invented 1974 the back-propagation having the ability to perform classification tasks beyond simple Perceptrons
 - Back-propagation was independently rediscovered in 1980s by David Rumelhat and David Parker

Artificial Neural Networks

Artificial Neural Networks

Input function of neuron



The input value **is multiplied by the weight before entering the neuron** these weighted inputs are then added together and generate the :

$$\sum_{i=1}^m O_{ij} w_{ij} \longrightarrow \text{input function}$$

O_i : Input of the neuron i from the pervious layer

Artificial Neural Networks

Artificial Neural Networks

Activation Function of neuron

- Biological neurons can have only **two activation states** (active, not active)
- Artificial neurons, however, can have **different activation states** range between (-1, 1) or (0,1)
- Lower bounds -1 or 0 stand for total inactivation
- Upper bound 1 stands for total activation

Example of usual activation functions
Sigmoid (logistic) function:

$$O_j = \frac{1}{1 + e^{-aZ_j}}$$

O_j is the output of the unit j of the actual layer. The parameter a is a real number and constant. Usually in ANN applications a is selected between 0.5 and 2

with $Z_j = \sum_{i=1}^m O_{ij} w_{ij} + \Theta_j$

where Θ_j is the bias

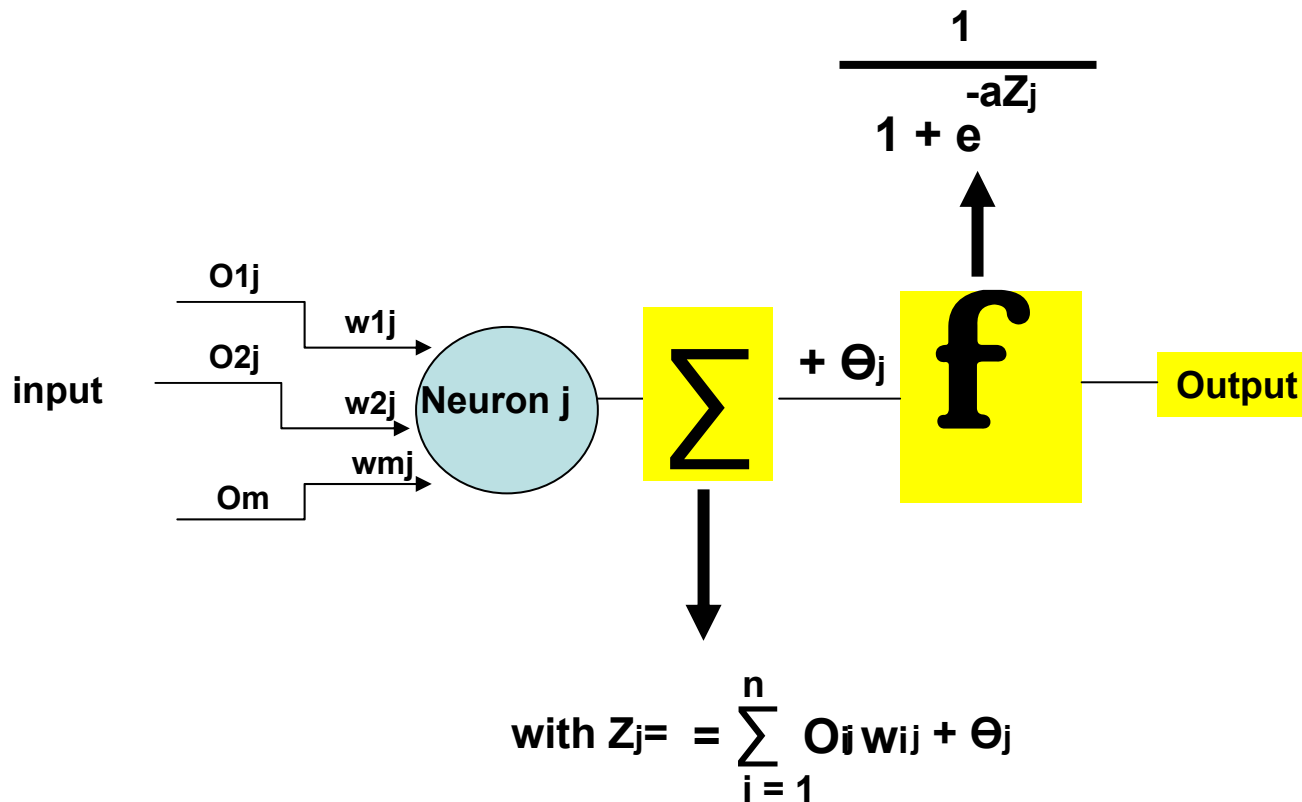
Output of the unit i of the previous layer

Net input

Artificial Neural Networks

Artificial Neural Networks

Output (function) of neuron



Artificial Neural Networks

Artificial Neural Networks

Activation Function of neuron

Why nonlinear activation functions are needed ?

To catch the nonlinearity aspects in the data, otherwise the hidden units could not make the ANN more powerful than just Simple *perceptrons* without any hidden units.

The nonlinearity makes representing of nonlinear functions possible and is the most powerful adjective of multilayer ANN

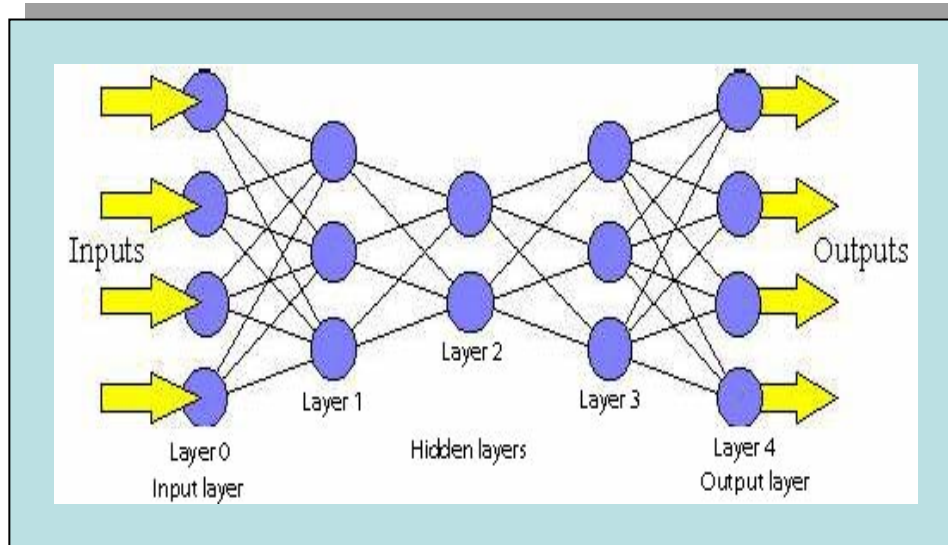
Artificial Neural Networks

Artificial Neural Networks

Architecture

Feed-Forward networks (FNN)

- FNN (known also as multi-layer *perceptrons*) are widely used models specially in practical applications
- They were the first and simplest type of ANN developed



In a FFN

- the information moves in only one direction
- information at a later level never backpropagates to the previous levels
- there are no cycles or loops in the network

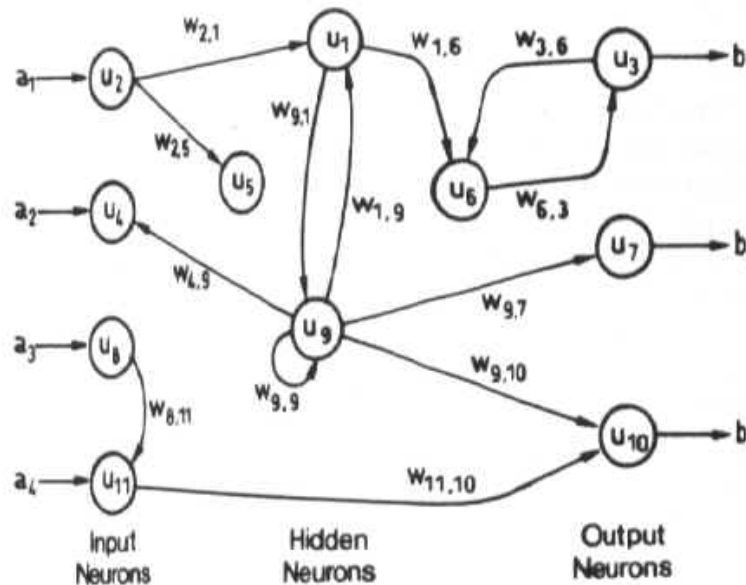
Artificial Neural Networks

Artificial Neural Networks

Architecture

Feedback networks (FBN)

The architecture of FBN (called also as interactive or recurrent Networks) is designed in a manner that they can send signals in both directions or in the loops



FBN are generally more difficult to train than FFN

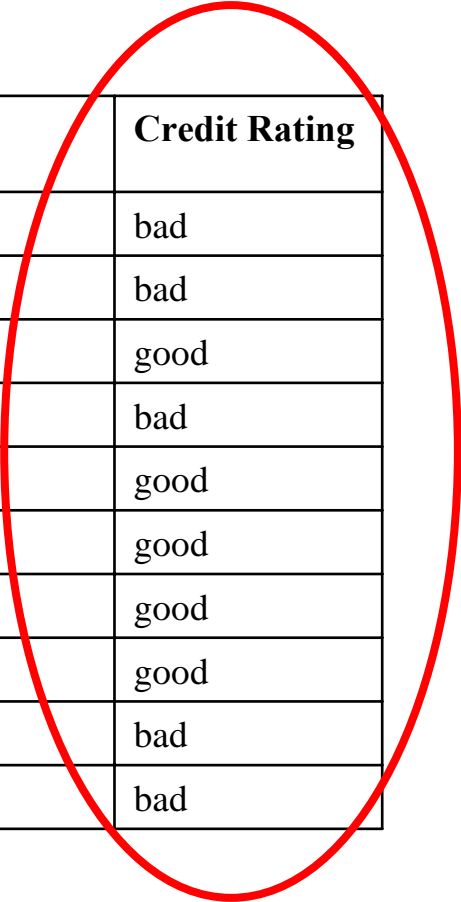
Artificial Neural Networks

Artificial Neural Networks

Learning Process

Supervised Learning

	Income >2000	Car	Gender	Credit Rating
1	no	yes	F	bad
2	no state	no	F	bad
3	no state	yes	M	good
4	no	yes	M	bad
5	yes	yes	M	good
6	yes	yes	F	good
7	no state	yes	F	good
8	yes	no	F	good
9	no state	no	M	bad
10	no	no	F	bad



Artificial Neural Networks

Artificial Neural Networks

Learning Process

Training set:

A set of examples used for learning, that is to fit the parameters [i.e., weights] of the classifier.

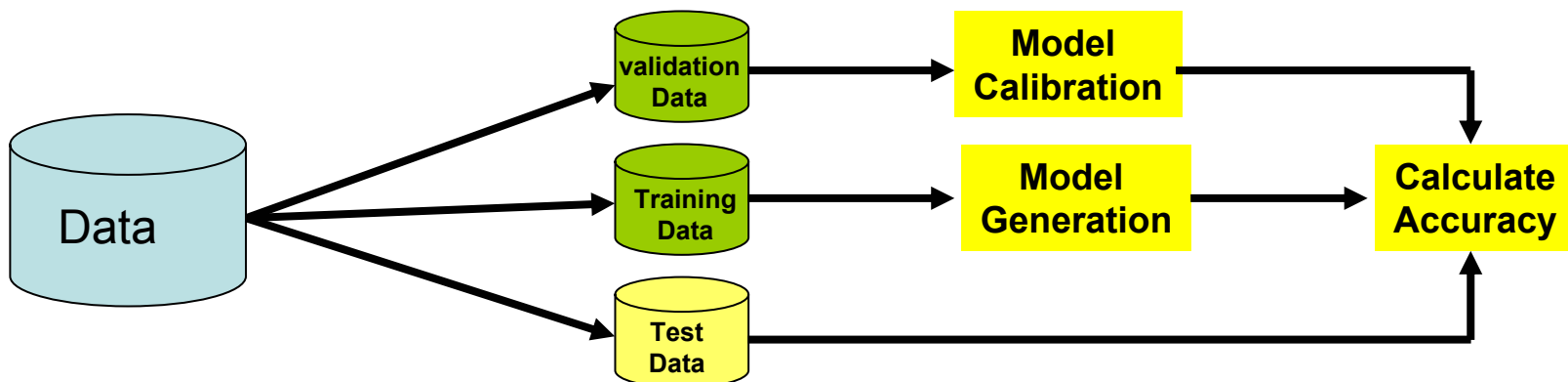
Ripley (1996), (p.354)

Validation set:

A set of examples used to tune the parameters [i.e., architecture, not weights] of a classifier, for example to choose the number of hidden units in a neural network.

Test set:

A set of examples used only to assess the performance [generalization] of a fully-specified classifier.



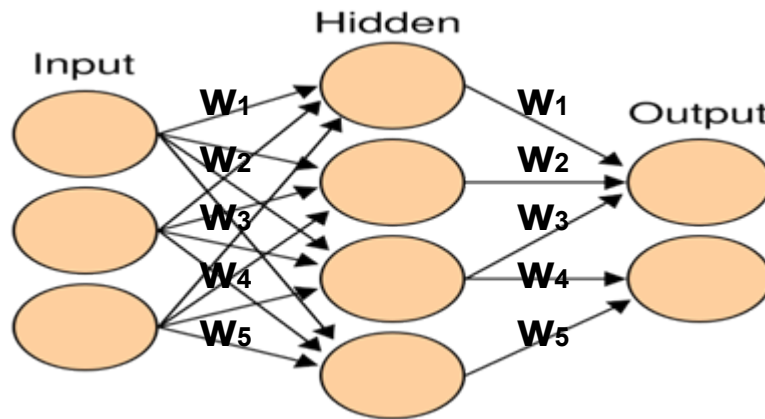
Artificial Neural Networks

Artificial Neural Networks

Learning Process

Learning Rule

How should we calculate the weight changing ?



new weight = old weight + weight change

$w(t+1) = w(t) + \Delta w(t)$ \longrightarrow Learning Rule

In the last years several Learning Rules have been invented

Artificial Neural Networks

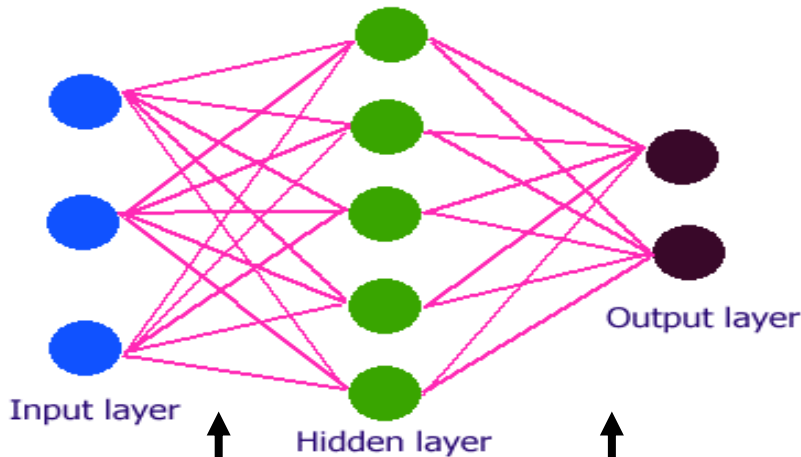
Artificial Neural Networks

Learning Process

Learning Rule

For each training Tuple:

The weights of the whole networks are modified in a manner to minimize the mean squared error (difference between the predicted and observed target value)



GD Learning Rule can be used to learn the weights of the outputs and hidden neurons

$$\text{Min: } E = \frac{1}{2} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

In most cases the output is a nonlinear function of weights and it is difficult to find the minimum of E analytically

Gradient descent method

$$w_j \rightarrow w_j - \beta * \partial E(w) / \partial w_j$$

β : Learning Rate between 0 and 1

Artificial Neural Networks

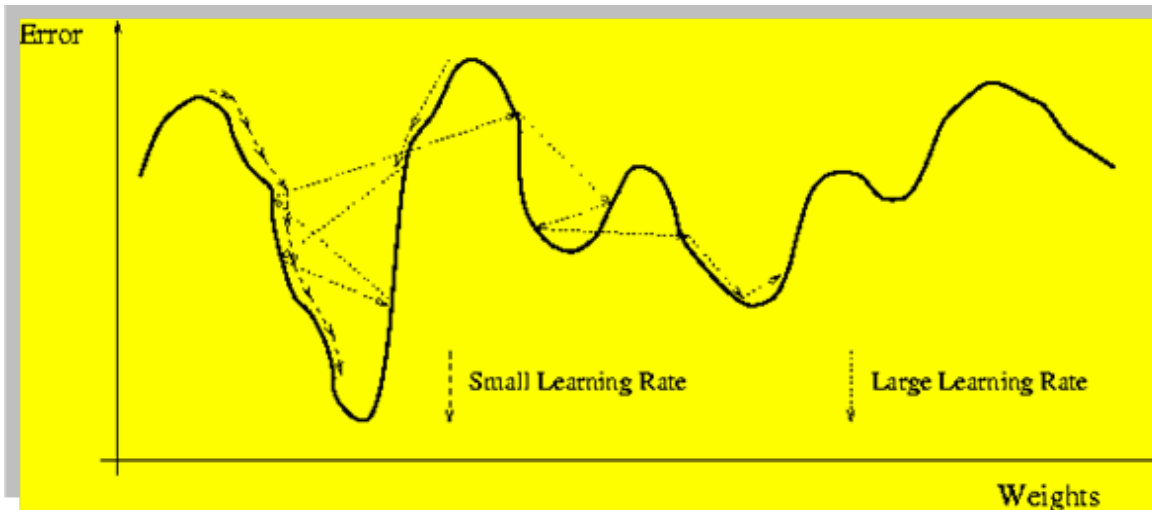
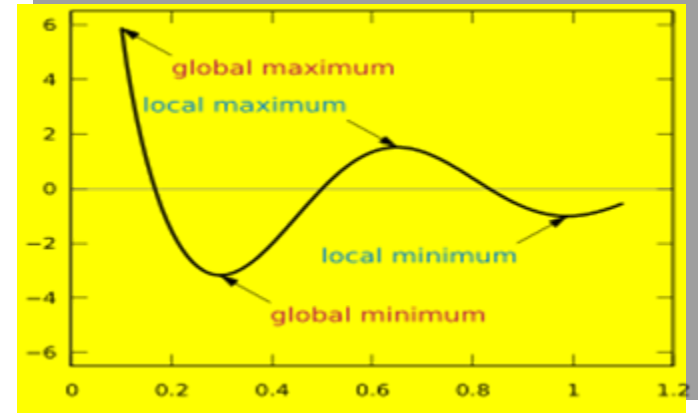
Artificial Neural Networks

Perceptron

Learning in Perceptron

Disadvantages of Gradient descent method:

- Getting Stuck in a local optima
- Small gradient \longrightarrow small change (slow)
- Big gradient \longrightarrow big change (cavort)



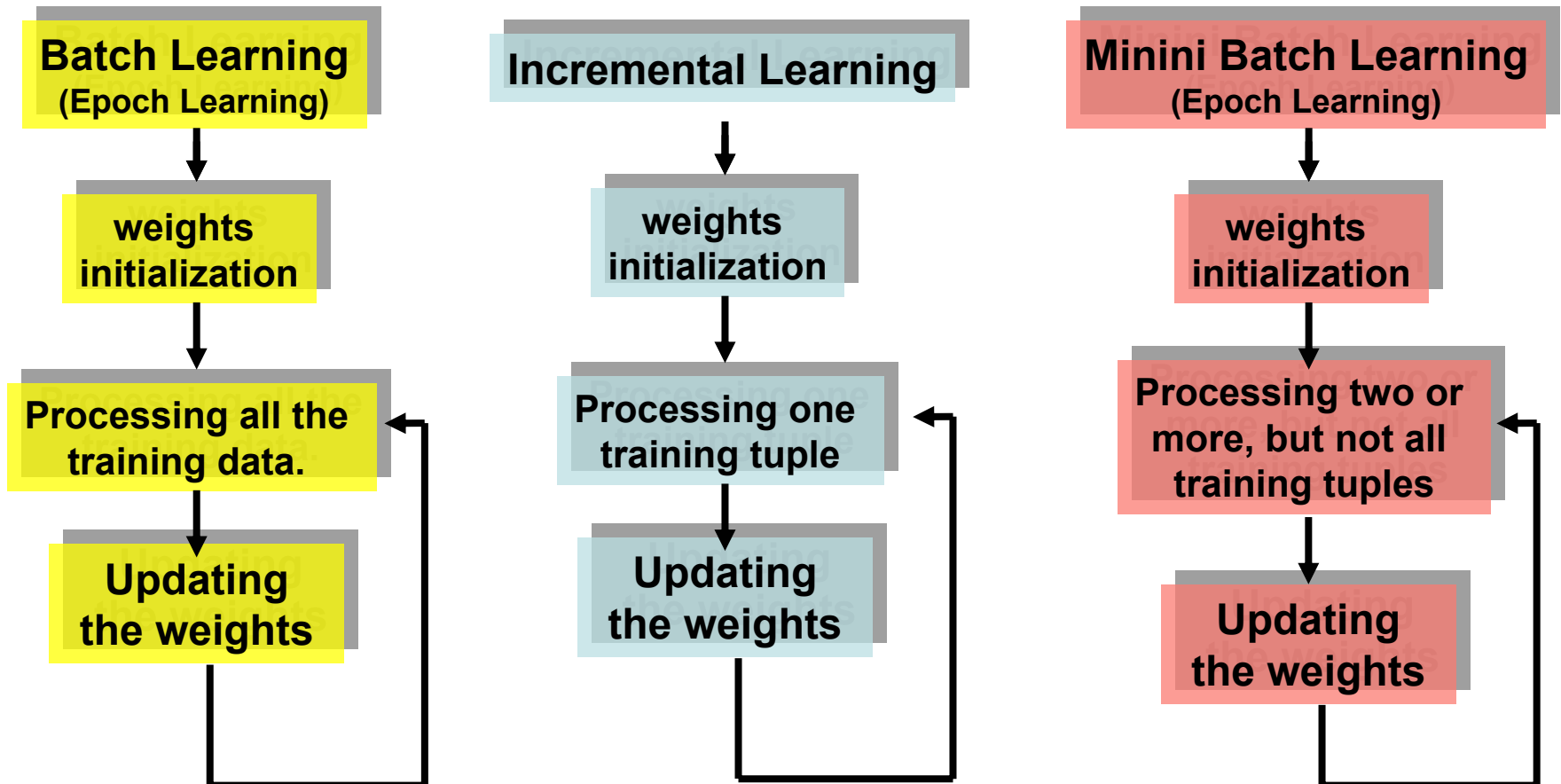
Artificial Neural Networks

Artificial Neural Networks

Learning Process

Learning Methods

Categorization of learning methods



Artificial Neural Networks

Artificial Neural Networks

Learning Process

Training Stop

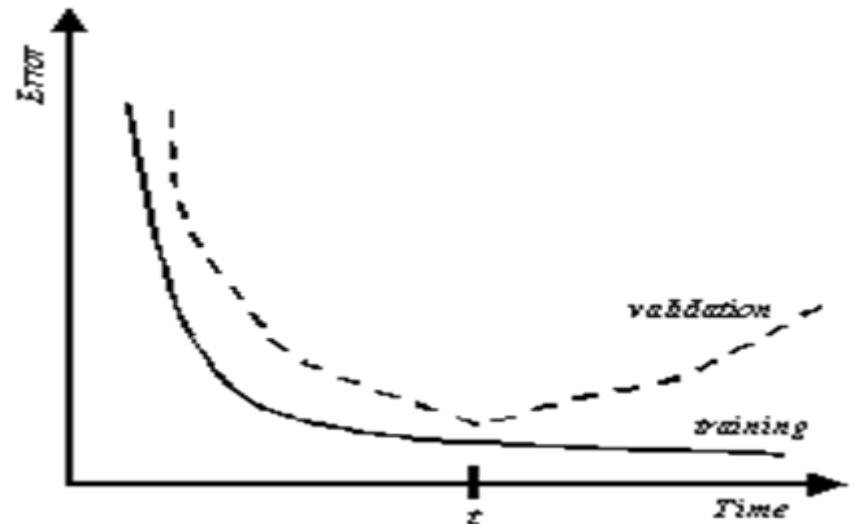
- Overfitting generates a relatively small error on the training dataset but larger error on new data
- To prevent Overfitting, it is sometime necessary to stop training by using a validation dataset

When stop training ?

- After a few epochs the ANN should be tested by using the validation dataset

Error on validation dataset increases higher than the last time error ?

Training Stop



Artificial Neural Networks

Artificial Neural Networks

Neuron's values

there are two types of neurons:

- **Binary neurons** can only takes values from the set $\{0,1\}$ or $\{1, -1\}$
- **Real-valued neurons** can take values in the intervals $[0,1]$ or $[1, -1]$

Artificial Neural Networks

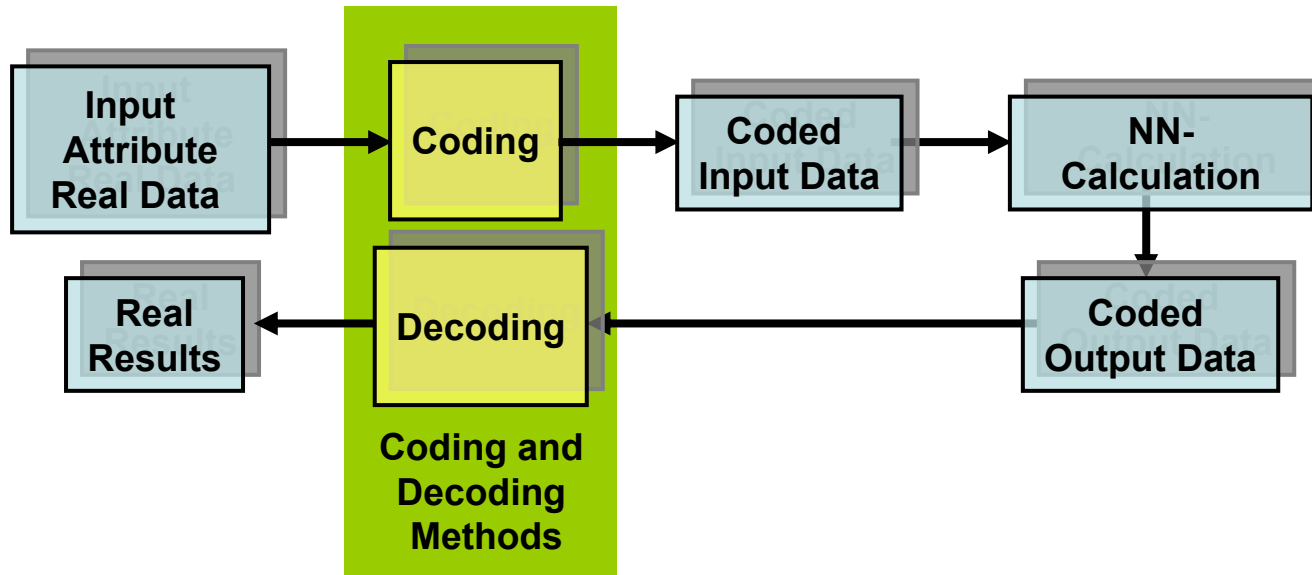
Artificial Neural Networks

Coding and decoding Methods

- **Coding:** transfer the continuous-valued or categorical data to the value domains:

$[0,1]$, $[-1, 1]$ or $\{0,1\}$, $\{-1, 1\}$

- **Decoding :** vice versa

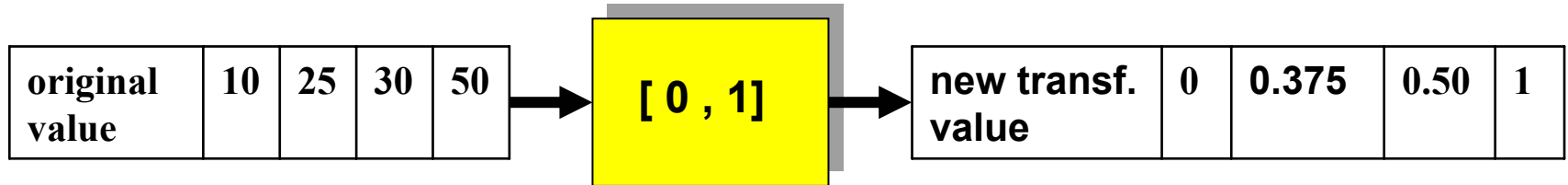


Artificial Neural Networks

Artificial Neural Networks

Coding and decoding Methods

Coding of continuous-valued data by using transformation function



X_n^{old} : original value ($n=1, 2, \dots, N$)

X_n^{new} : new, transferred value ($n=1, 2, \dots, N$)

$X_{\text{max}}^{\text{old}}$: maximal original value

$X_{\text{min}}^{\text{old}}$: minimal original value

$$X_n^{\text{new}} = A * X_n^{\text{old}} + B$$

$$A = \frac{X_{\text{max}}^{\text{new}} - X_{\text{min}}^{\text{new}}}{X_{\text{max}}^{\text{old}} - X_{\text{min}}^{\text{old}}}$$

$$B = X_{\text{min}}^{\text{new}} - A * X_{\text{min}}^{\text{old}}$$

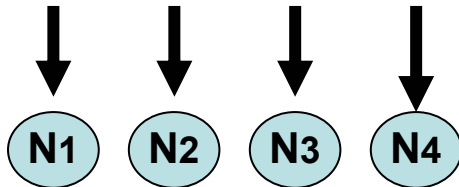
Artificial Neural Networks

Artificial Neural Networks

Coding and decoding Methods

Coding of nominal attributes

Marital status				
single	1	0	0	0
married	0	1	0	0
widowed	0	0	1	0
divorced	0	0	0	1



For each value of a nominal attribute an input neuron is necessary

Many nominal attributes with several values lead to a large number of input neurons

Artificial Neural Networks

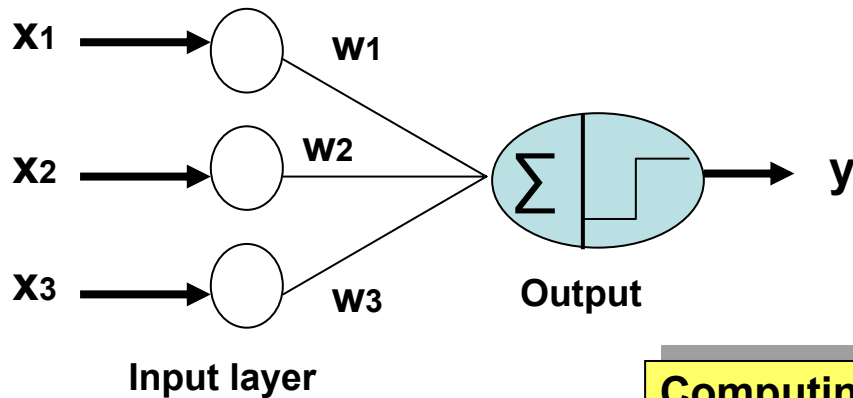
Artificial Neural Networks

Network type

Perceptron

Perceptron:

A simple ANN architecture consists only of input and output nodes (no hidden layer):



Computing the output y in Perceptron:

$$\hat{y} = \begin{cases} 1, & \text{if } w_1 x_1 + w_2 x_2 + w_3 x_3 - t > 0 \\ -1, & \text{if } w_1 x_1 + w_2 x_2 + w_3 x_3 - t \leq 0 \end{cases}$$

t : Bias factor

Generally:

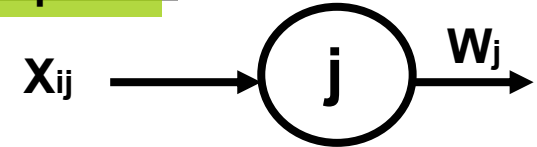
$$\hat{y} = \text{sign} (w_1 x_1 + w_2 x_2 + \dots w_j x_j + \dots w_m x_m - t)$$

Artificial Neural Networks

Artificial Neural Networks

Perceptron

Learning in Perceptron



x_{11}	x_{12}	x_{1j}	x_{1m}	y_1
.....						
x_{i1}	x_{i2}	x_{ij}	x_{im}	y_i
.....						
x_{n1}	x_{n2}	x_{nj}	x_{nm}	y_n

$$w_j^{(k+1)} = w_j^{(k)} + \beta (y_i - \hat{y}_i^{(k)}) x_{ij}$$

- $w_j^{(k+1)}$: new weight in iteration $k + 1$ of neuron j
- $w_j^{(k)}$: old weight in iteration k of neuron j
- β : learning rate
- y_i : observed output of the tuple i
- $\hat{y}_i^{(k)}$: calculated output of the tuple i in iteration k
- x_{ij} : value of attribute j of the tuple i

The above Learning Rule is based on Gradient Descent method by minimizing

$$\text{Min: } E = \frac{1}{2} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Artificial Neural Networks

Artificial Neural Networks

Perceptron

Learning in Perceptron

$$\text{Min: } E = \frac{1}{2} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad \text{for tuple } i \quad \text{Min: } E = \frac{1}{2} (y_i - \hat{y}_i)^2$$

$$\hat{y}_i = \sum_{j=1}^m w_{ij} x_{ij} \quad \partial \hat{y}_i / \partial w_{ij} = x_{ij}$$

$$\partial E / \partial w_j = - (y_i - \hat{y}_i) * \partial \hat{y}_i / \partial w_j = - (y_i - \hat{y}_i) x_{ij} \quad (1)$$

$$w_j \rightarrow w_j - \beta * \partial E(w) / \partial w_j \quad (2)$$

(1) and (2) lead to

$$w_j^{(k+1)} = w_j^{(k)} + \beta (y_i - \hat{y}_i^{(k)}) x_{ij}$$

Artificial Neural Networks

Artificial Neural Networks

Perceptron

Learning in Perceptron

The relation $w_1 x_1 + w_2 x_2 + \dots + w_j x_j + \dots + w_m x_m$ is linear in w and x

For linearly separable classification problems, the learning algorithm

$$w_j^{(k+1)} = w_j^{(k)} + \beta (y_i - \hat{y}_i^{(k)}) x_{ij}$$

converges if the learning rate β is sufficiently small

If the classes are not linearly separable, the above learning rule does not converge .

For such tasks ANN with hidden layers are necessary. Backpropagation is such an alternative.

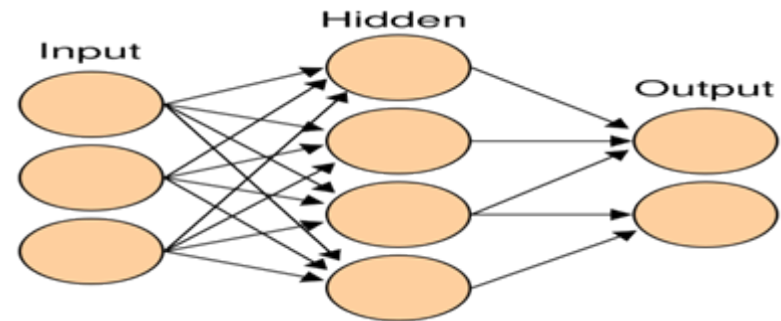
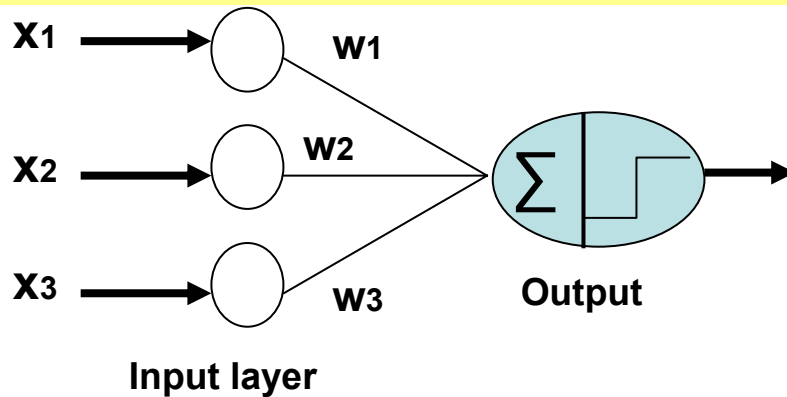
Artificial Neural Networks

Artificial Neural Networks

Backpropagation Algorithms

Introduction

Backpropagation (backprop , BP) is the mostly used multilayer feed forward ANN, specially in the praxis . As we have seen, Perceptron can not handle the nonlinearly separable classes, but, BP can



The modification of the weights going to a hidden unit can not be performed by the method used in Perceptron because – in contrast to the output units – we have no information about the observed outputs of the hidden units. It means, it is not possible to determine the error term related to each hidden unit.

Solution: propagation of the errors of the output units backwards → Backpropagation

Artificial Neural Networks

Artificial Neural Networks

Backpropagation Algorithms

Learning Steps

Weight initialization : network weights are initialized randomly, usually from the interval $[-1, 1]$ or $[-0.5, 0.5]$

Calculation of the net input: As we have seen before the net Input of the neuron j is

$$z_j = \sum_{i=1}^p X_{ij} w_{ij} + \Theta_j$$

(1)

Calculation of the output function: Using a logistic function we have got

$$O_j = \frac{1}{1 + e^{-az_j}}$$

(2)

This function is differentiable and nonlinear

Repeat (1) and (2) until we the output layer is reached and included

Artificial Neural Networks

Artificial Neural Networks

Backpropagation Algorithms

Learning Steps

Backpropagate the error

E: Error of the unit j

δ_j
↑

$$\partial E / \partial w_{ij} = \partial E / \partial z_j * \partial z_j / \partial w_{ij} = x_{ij} * \partial E / \partial z_j = x_{ij} \delta_j$$

A- Neuron j belongs to the **Output Layer**

Contribution to **E** by j : $E = 1/2 (T_j - O_j)^2$ T_j : target Value of the neuron j

$\delta_j = \partial / \partial z_j E = -(T_j - O_j) \partial O_j / \partial z_j$ and for **a = 1**

$$\delta_j = -(T_j - O_j) (1 - O_j) O_j$$

Therefore : regarding the Gradient Descent Learning Rule : $W_{ij} \rightarrow W_{ij} + \beta \delta_j X_{ij}$

Artificial Neural Networks

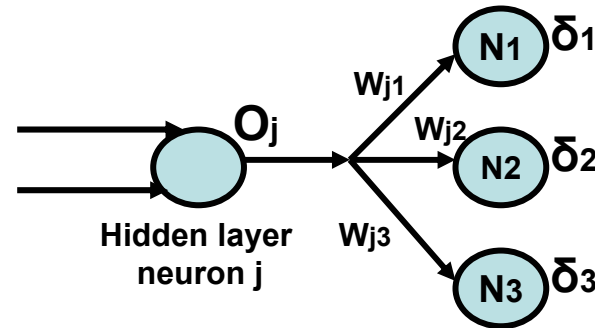
Artificial Neural Networks

Backpropagation Algorithms

Learning Steps

Backpropagate the error

Case B- Neuron j belongs to a *Hidden Layer*



$$E_j = \delta_j = O_j (1 - O_j) \sum_k \delta_k W_{jk}$$

- O_j : Oput of the neuron j
- δ_k : Error of the neuron k in the next layer
- W_{jk} : weight of the connection from the neuron j to a neuron k in the next layer

Gradient Descent Learning Rule : $W_{ij} \rightarrow W_{ij} + \beta \delta_j X_{ij}$

Bias updating for Neuron J : $\theta_j = \theta_j + \beta E_j$

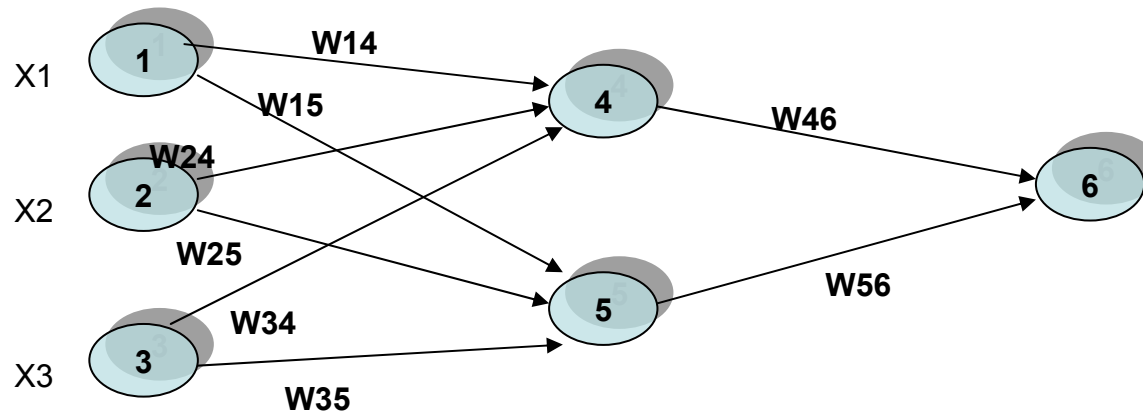
Artificial Neural Networks

Artificial Neural Networks

Backpropagation Algorithms

Example

(source: Han et al (2006))



x1	x2	x3	w14	w15	w24	w25	w34	w35	w46	w56	Θ4	Θ5	Θ6
1	0	1	0.2	-0.3	0.4	0.1	-0.5	0.2	-0.3	-0.2	-0.4	0.2	0.1

Unit	Net Input, Z_j	Output, O_j
4	$0.2 + 0 - 0.5 - 0.4 = -0.7$	$\frac{0.7}{1 + e^{-0.7}} = 0.332$
5	$-0.3 + 0 + 0.2 + 0.2 = 0.1$	$\frac{-0.1}{1 + e^{-0.1}} = 0.525$
6	$(-0.3)(0.332) - (0.2)(0.525) + 0.1 = -0.105$	$\frac{0.105}{1 + e^{-0.105}} = 0.474$

Artificial Neural Networks

Artificial Neural Networks

Backpropagation Algorithms

Example

(source: Han et al (2006))

Unit j	Error j
6	$(0.474) (1 - 0.474) (1 - 0.474) = 0.1311$
5	$(0.525) (1 - 0.525) (0.1311) (-0.2) = -0.0065$
4	$(0.332) (1 - 0.332) (0.1311) (-0.3) = -0.0087$

Weight or bias	New value
W46	$-0.3 + (0.9) (0.1311) (0.332) = -0.261$
W56	$-0.2 + (0.9) (0.1311) (0.525) = -0.138$
W14	$0.2 + (0.9) (-0.0087) (1) = 0.192$
W15	$-0.3 + (0.9) (-0.0065) (1) = -0.306$
W24	$0.4 + (0.9) (-0.0087) (0) = 0.4$
W25	$0.1 + (0.9) (-0.0065) (0) = 0.1$
W34	$-0.5 + (0.9) (-0.0087) (1) = -0.508$
W35	$0.2 + (0.9) (-0.0065) (1) = 0.194$
Θ_6	$0.1 + (0.9) (0.1311) = 0.218$
Θ_5	$0.2 + (0.9) (-0.0065) = 0.194$
Θ_4	$-0.4 + (0.9) (-0.0087) = -0.408$

Weakness and Strength of ANN

Based on: <http://www-sal.cs.uiuc.edu/~hanj/bk2/>

- **Strength**

- High tolerance to noisy data
- Well-suited for continuous-valued inputs and outputs
- Successful on a wide array of real-world data
- Techniques have recently been developed for the extraction of rules from trained neural networks

- **Weakness**

- Long training time, specially for the datasets with a large number of categorical attributes
- Require a number of parameters typically best determined empirically, e.g., the network topology or "structure."
- Poor interpretability: Difficult to interpret the symbolic meaning behind the learned weights and of "hidden units" in the network

- **Introduction**
- **Rise and fall of Neural Networks**
- **Input function of neuron**
- **Activation function of neuron**
- **Output (function) of neuron**
- **Feed-Forward networks**
- **Feedback networks**
- **Learning process**
- **Coding and decoding methods**
- **Perceptron**
- **Backpropagation**
- **Weakness and Strength of ANN**

