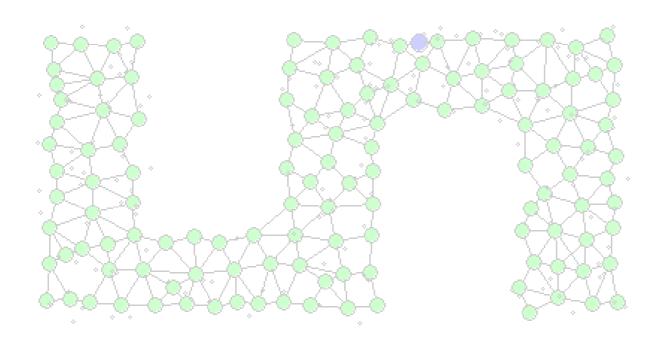
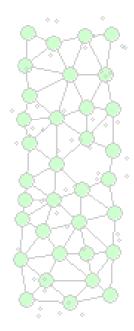
# Kapitel VI Neuronale Netze

(basierend auf Material von Andreas Hotho)





#### Agenda

- 1. Einführung & Grundbegriffe
  - Motivation & Definition
  - Vorbild Biologie
  - Historie der NN
  - Überblick über verschiedene Netzwerktypen
- 2. Einfaches Perzeptron
- 3. Multi Layer Perzeptron
- 4. Beispiele
- 5. Rekurrente Netze
- 6. Entwurfsüberlegungen

#### Was sind künstliche Neuronale Netze

#### Künstliche Neuronale Netze sind

- massiv parallel verbundene Netzwerke aus
- · einfachen (üblicherweise adaptiven) Elementen in
- · hierarchischer Anordnung oder Organisation,

die mit der Welt in der selben Art wie biologische Nervensysteme interagieren sollen.

(Kohonen 84)

#### Wofür nutzt man künstliche Neuronale Netze

- Forschung:
  - Modellierung & Simulation biologischer neuronaler Netze
  - Funktionsapproximation
  - Speicherung von Informationen
  - •

- · Anwendungen (z.B.):
  - Interpretation von Sensordaten
  - · Prozessteuerung
  - Medizin
  - Elektronische Nase
  - Schrifterkennung
  - · Risikomanagement
  - · Zeitreihenanalyse und -prognose
  - Robotersteuerung

#### **Kapitel VI.3: Neuronale Netze**

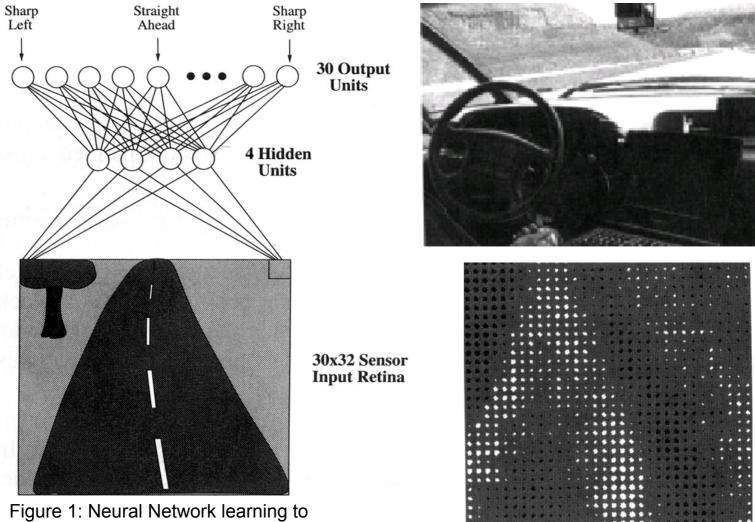


Figure 1: Neural Network learning to steer an autonomous vehicle (Mitchell 1997)

#### **Kapitel VI.3: Neuronale Netze**

- <u>charakteristische</u> Eigenschaften von Problemen, die mit BACKPROPAGATION ANNs gelöst werden können:
  - <u>Trainingsbeispiele</u> sind repräsentiert durch <u>Attribut-Wert-Paare</u>
    - Werte können reellwertig sein
  - + zu generierende Funktion (target function) kann ...
    - · <u>diskrete</u> Funktionswerte
    - reellwertige Funktionswerte oder
    - Vektor solcher Funktionswerte

#### haben

- + Trainingsbeispiele dürfen <u>fehlerhaft</u> sein
- <u>lange</u> Trainingszeiten sind akzeptabel
- + <u>schnelle</u> Berechnung der Funktionswerte der gelernten Funktion kann erforderlich sein
- für Menschen <u>verständliche</u> Interpretation der gelernten Funktion ist <u>nicht</u> wichtig ("Black Box - Ansatz")

#### Arbeitsweise Neuronale Netze

## gekennzeichnet durch:

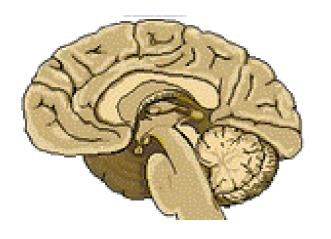
- massiv parallele Informationsverarbeitung
- Propagierung der Informationen über Verbindungsstellen (Synapsen)
- verteilte Informationsspeicherung
- Black Box Charakter

#### es werden die:

- Aufbauphase (Topologie),
- · Trainingsphase (Lernen) und
- Arbeitsphase (Propagation) unterschieden (Phasen können auch überlappen)

## Vorbild Biologie

#### Gehirn



In der Vergangenheit wurden Forscher aus den verschiedensten Fachgebieten durch das biologische Vorbild motiviert.

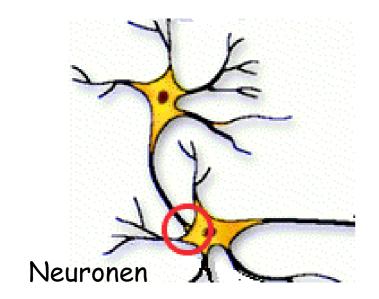
# Geflecht aus Neuronen

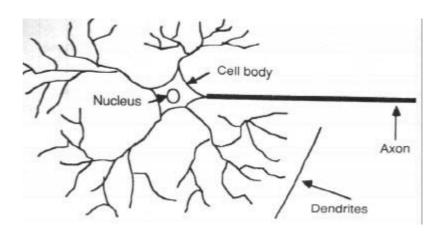
#### Das Gehirn besteht aus

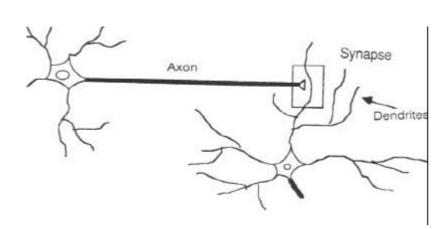
- · ca. 1011 Neuronen, die mit
- · ca. 10<sup>4</sup> anderen Neuronen
- ca. 10<sup>13</sup> Synapsen verschaltet sind.

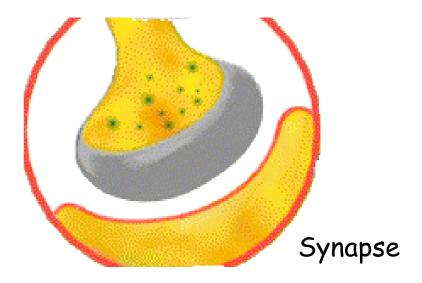


# Vorbild Biologie



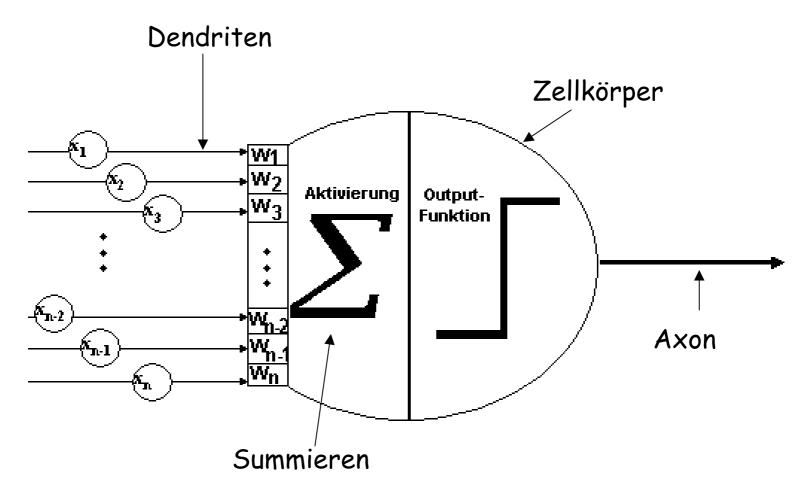






# Vorbild Biologie

#### Vom natürlichen zum künstlichen Neuronalen Netz



#### Elemente eines Neuronalen Netzes

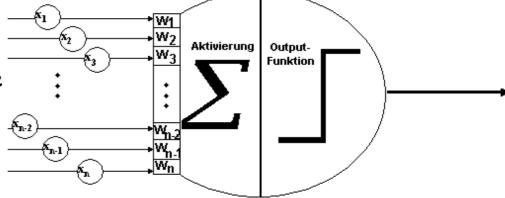
- Menge von Neuronen (Berechnungseinheiten, oder elementare Prozessoren)
- Verbindungen zwischen den Neuronen im Netz (Synapsen)
- · Gewichte der Verbindungen
- Propagierungsregel für das Weiterleiten des Inputs
- · Aktivierungs-Zustände der Neuronen
- Output-Funktionen und Schwellwerte der Neuronen

· Aktivierungsregel (Verknüpfung des Inputs und des Zustands zum

Output)

· Lernregel

 Umgebung, in der das lernende System arbeiten soll



Biologische Vorbilder wurden in den letzten 150 Jahren erforscht

- 1850 Nervenimpulse als elektrisches Phänomen
- · 1901 Richtung der Informationsübertragung
- · 1910-1930 synaptische Informationsübertragung
- 1949-1956 Informationsweiterleitung (ionische Depolarisation) geklärt
- ca. 1979 Abbildung der sensorischen Informationen als Karte im Gehirn (Verteilung über die Großhirnrinde)

Neben den Biologen haben auch Psychologen und Ingenieure viel zur Weiterentwicklung des Verständnisses über Neuronale Netze beigetragen.

1943	McCulloch und Pitts entwickelten erste künstliche
	Neuronale Netze auf ihrem Verständnis der Neurologie
1949	Hebb'sche Lernregel
1957	entwickelt Rosenblatt das Perceptron
1969	Minski und Papert führten eine rigorose Analyse des
	Perceptrons durch
1972	Kohonen Assoziativspeicher
1974	Werbos entwickelt Backpropagation Algorithm
1976	Grossberg entwickelt ART1 Modell
1982	Kohonen's selbstorganisierende Karten
1985	Hopfield, Tank realisierten neue Ideen wie
	rückgekoppelte Netzwerke mit selbstorganisierendem
	Lernen
1986	Rumelhart und Hinton führten die alten Ideen des
	Perceptrons in mehrschichtigen Modellen weiter
	(Backpropagation Algorithmus) MLP

#### McCulloch und Pitts (1943) Netze

- · binare Ein- und Ausgabe
- n erregende Leitungen  $(x_1, ..., x_n)$
- m hemmende Leitungen  $(y_1, ..., y_m)$
- · Schwellenwert s
- Ausgabefunktion
- · darf zyklisch sein
- · Ausgabefunktion ist eine Treppenfunktion
- · Berechnung benötigt eine Zeiteinheit

Falls  $m \ge 1$  und eines der Signale  $y_1, ..., y_m$  gleich eins ist

dann: Output = 0

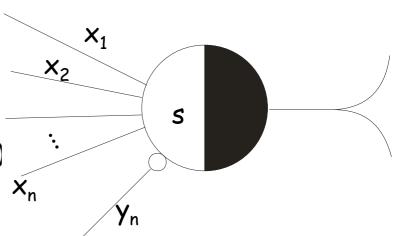
sonst: Summe  $(x_1, ..., x_m)$  wird berechnet

und mit Schwellenwert s verglichen

Falls Erregung ≥ s ist

dann: Output = 1

sonst: Output = 0



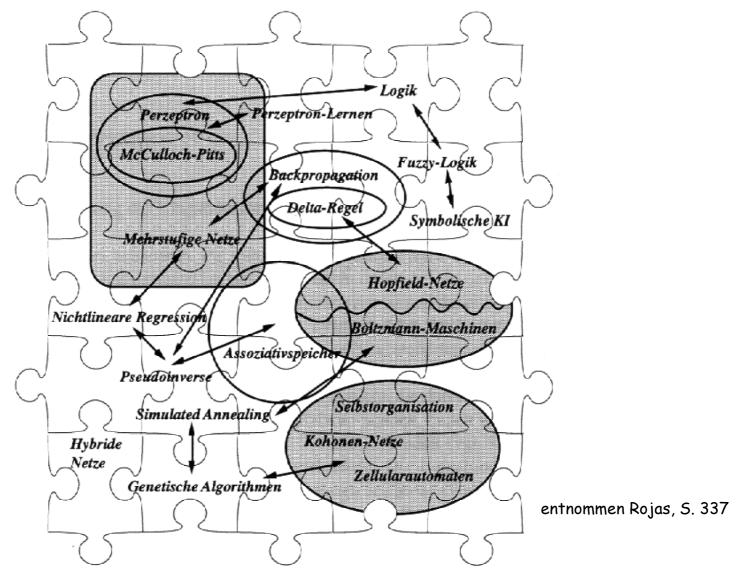
#### Vorteile

- erlaubt die Konstruktion von beliebigen logischen Funktionen
- jeder endliche Automat kann mit einem Netz aus McCulloch-Pitts Zellen simuliert werden

#### Nachteil

· Lernen sehr aufwendig (Topologieänderung nötig)

# Typologie künstlicher neuronaler Netzwerk



# Typologie künstlicher neuronaler Netzwerk

gewichtete <---> ungewichtete Netze

synchrone <---> asynchrone Netze

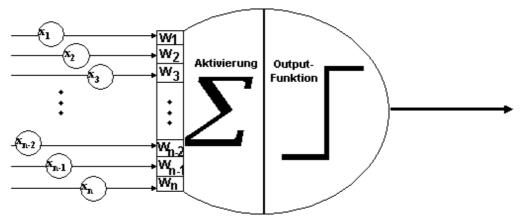
mit <---> ohne gespeichertem Zustand

mit <---> ohne Rückkopplung

in Schichten <---> ohne Schichten

binäre <---> reelle Zustände

# Aktivierungsfunktionen



Es gibt i.w. vier Arten der Aktivierung von Neuronen:

Allen gemeinsam ist die **Gewichtung** der Inputs. Die Inputs  $x_1, x_2, ..., x_n$  werden stets mit ihren (Synapsen-)Gewichten  $w_1, w_2, ..., w_n$  multipliziert:  $a_1 = x_1^* w_1$ ,  $a_2 = x_2^* w_2, ..., a_n = x_n^* w_n$ .

1. Die am häufigsten angewandte Regel ist die (Skalarprodukt-Regel)

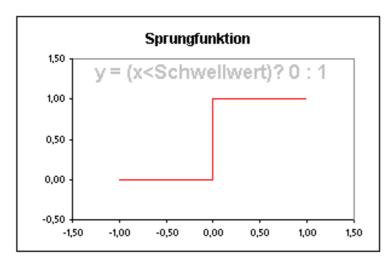
Die gewichteten Inputs  $a_1,a_2,...,a_n$  werden zur Aktivität des Neurons aufaddiert:

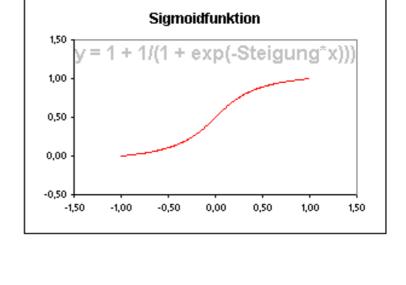
$$a = a_1 + a_2 + ... + a_n$$

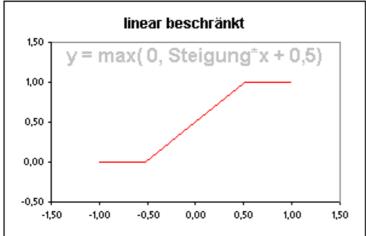
# Aktivierungsfunktionen

- 2. Sehr häufig ist ebenfalls die Winner-take-all-Regel: bei der die Aktivität a zunächst nach der Skalarproduktregel ermittelt wird, dann aber mit allen Aktivitäten in derselben Schicht verglichen wird und auf 0 herabgesetzt wird, wenn ein anderes Neuron höhere Aktivität hat.
- 3. Bisweilen kommt auch die **sigma-pi-Regel** vor: Die Inputs werden in k Gruppen aufgeteilt und die modifizierten Inputs jeder Gruppe miteinander multipliziert, danach werden die Produkte  $p_1,p_2,...,p_k$  zur Aktivität des Neurons aufaddiert  $a = p_1+p_2+...+p_k$ . Der Effekt ist, daß ein Input von fast 0 die gesamte Gruppe blockiert.
- 4. Selten ist die Regel, bei der ebenfalls die Inputs in Gruppen eingeteilt werden (z.B. 2 Gruppen: inhibitorische und excitatorische) und gruppenweise nach Skalarproduktregel zusammengefasst werden. Die Ergebnisse werden dann mit einer nichtlinearen Funktion miteinander verknüpft (z.B. a=(e-i)/|i|)

# Outputfunktionen







Alle Funktionen könnten wie die Sprungfunktion mit einem Schwellwert verschoben werden. Wie wir aber sehen werden kann auf diese Verschiebung verzichtet werden.

--> Demo (funktion.html)

## Lernparadigmen

## Lernverfahren lassen sich grob unterteilen in:

- supervised (überwachtes)
   Der Trainingsdatensatz enthält nicht nur den Input für das Neuronale Netz sondern auch den gewünschten Output.
   Anhand dieses Outputs kann der Fehler des Netzes bestimmt und korrigiert werden.
- unsupervised (unüberwachtes)
   Neuronales Netz muß anhand des präsentierten Inputs die Muster in den Daten erkennen.

#### Overfitting

- Mit Overfitting beschreibt man das Problem der Überanpassung eines Modells an einen Datensatz. Das Modell paßt sich sehr gut an den kleinen Weltausschnitt an, kann aber wegen der fehlenden Generalisierung nur schlecht auf neue Situationen reagieren.
- Kontrolle der Generalisierung während des Trainings durch Teilen des Datensatzes
  - Trainingset
  - ·Validierungsset
  - ·Testset

## **Fuzzy**

- Neuronale Netze sind in der Lage, Regelaufgaben zu erlernen.
- · Aber es ist keine Interpretation möglich
- kombiniert man Fuzzy und Neuronale Netze, so ist man in der Lage, die gelernten Regel aus dem Neuronalen Netz zu extrahieren
- Regel können dann auch interpretiert und von Experten geprüft werden

mehr dazu unter: http://fuzzy.cs.uni-magdeburg.de/index.html

## Genetische Algorithmen

- Genetische Algorithmen können für das Training von Neuronalen Netzen eingesetzt werden und sind in der Lage, das globale Optimum zu finden.
- · eine ganze Population von Netzen wird konstruiert
- beste Netze erhalten dann auch die besten Chancen
- · schlechte Netze sterben aus

mehr dazu unter:

http://www-ra.informatik.uni-tuebingen.de/forschung/welcome.html

#### Statistik

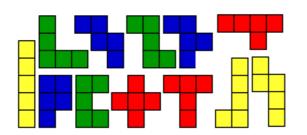
In welchem Zusammenhang stehen Neuronale Netze und statistische Analyse?

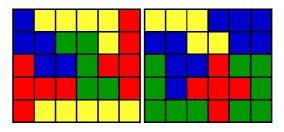
- Teilgebiete der Statistik beschäftigen sich mit Datenanalyse
- · aus Sicht der Neuronalen Netz bedeutet dies:
  - --> statistisches Schlußfolgern heißt für NN, Verallgemeinerungen aus verrauschten Daten zu lernen
- · Beispiel:
  - vorwärtsgerichtete Netze ohne versteckte Schicht entsprechen linearen Modellen
  - · Kohonen Netz ist dem K-Means Clusterverfahren ähnlich
  - · Hebbian Lernen ist Hauptkomponentenanalyse ähnlich

#### Vor- und Nachteile Neuronaler Netze

#### Vorteile

sehr gute Mustererkenner





- verarbeiten verrauschte, unvollständige und widersprüchliche Inputs
- verarbeiten multisensorischen Input (Zahlen, Farben, Töne, ...)
- erzeugen implizites Modell für Eingaben (ohne Hypothesen des Anwenders)
- · fehlertolerant auch gegenüber Hardwarefehlern
- · leicht zu handhaben

#### Vor- und Nachteile Neuronaler Netze

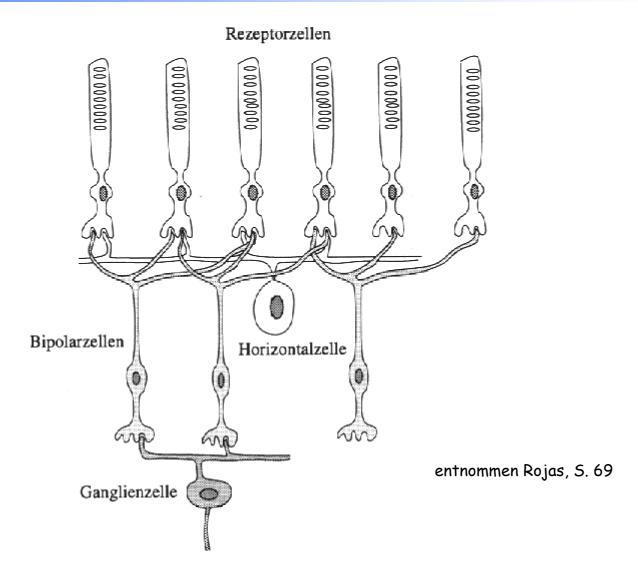
#### Nachteile

- lange Trainingszeiten
- · Lernerfolg kann nicht garantiert werden
- Generalisierungsfähigkeit kann nicht garantiert werden (Overfitting)
- keine Möglichkeit, die Begründung einer Antwort zu erhalten (Blackbox)

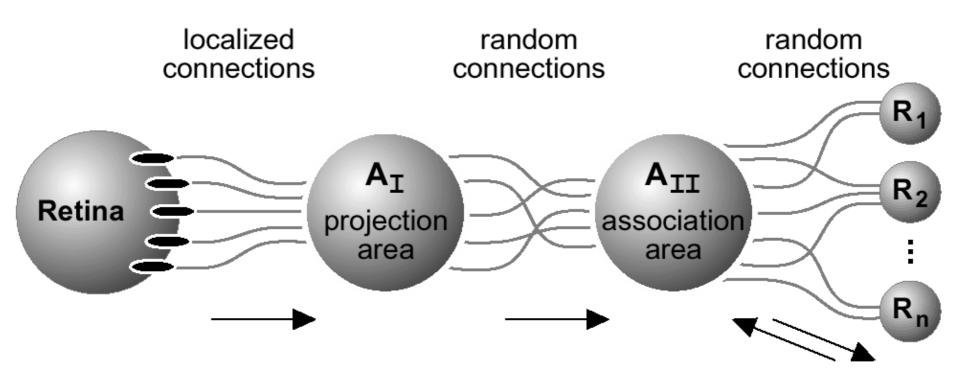
## Agenda

- 1. Einführung
- 2. Einfaches Perzeptron
  - Motivation
  - Definition Perzeptron
  - Geometrische Interpretation
  - Lernen Delta Regel
  - XOR Problem
- 3. Multi Layer Perzeptron
- 4. Beispiel
- 5. Rekurrente Netze
- 6. Entwurfsüberlegungen

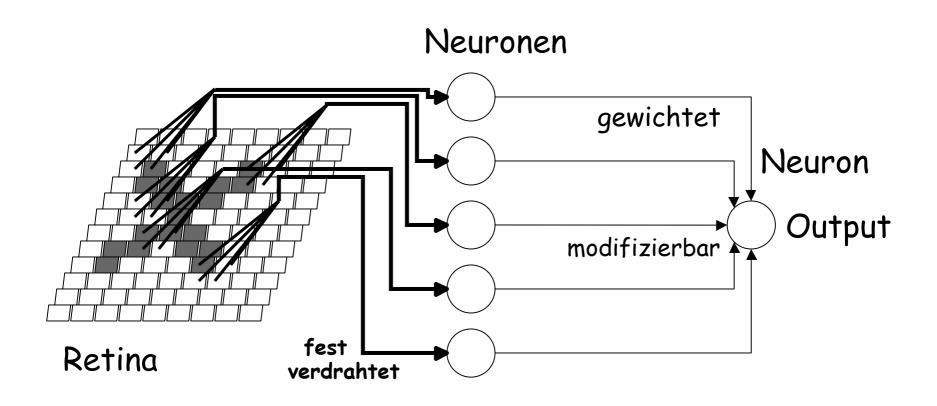
# Verschaltung der Netzhaut



# Perzeptron von Rosenblatt 1958



# Perzeptron von Rosenblatt 1958



Inputschicht

Outputschicht

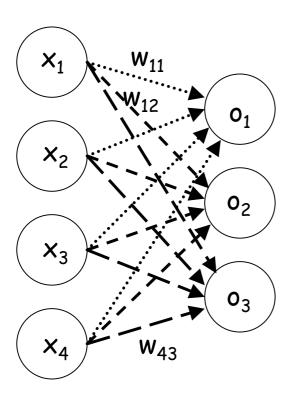
#### Perzeptron von Rosenblatt 1958

- · feste Verbindung zwischen Retina und Inputschicht
- Vorverarbeitung erfolgt zwischen Retina und Inputschicht
- Ausgabe der Inputschicht ist 0 oder 1
- Art der Vorverarbeitung (Verbindung + Neuronen) wurde im Originalmodell nicht spezifiziert
- Lernvorgang konzentriert sich auf die gewichteten Verbindungen zwischen Input- und Outputschicht

## Perzeptron

· jedes Outputneuron hat einen eigenen unabhängigen Netzbereich





- · d.h., für weitere Betrachtungen genügt ein Netz mit einem Neuron in der Outputschicht
- Input  $x = (x_1, ..., x_n)$  Gewichte  $w = (w_1, ..., w_n)$ Input

$$x = (x_1, ..., x_n)$$

$$w = (w_1, ..., w_n)$$

Output

Inputschicht Outputschicht

#### Perzeptron

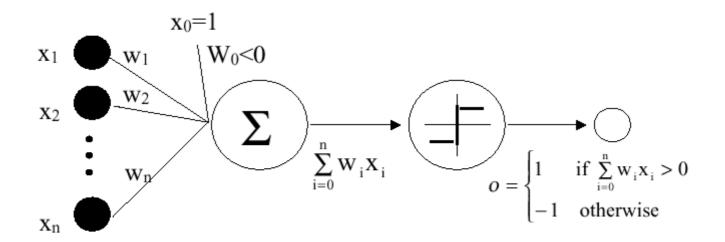


Figure 2: A Perceptron (Mitchell 1997)

• Input  $x = (x_1, ..., x_n)$ • Gewichte  $w = (w_1, ..., w_n)$ • Output (o)

## Perzeptron

- Outputneuron hat:
  - Schwellenwert s
  - Aktivität a := xw :=  $x_1w_1 + ... + x_nw_n$ (Skalarproduktregel)
  - verwendet Sprungfunktion
- · Output berechnet sich wie folgt:

```
o = 0, falls a < s
o = 1, sonst.
```

• äquivalente Beschreibung: erweitert man die Vektoren x und w zu  $y = (x_1, ..., x_{n,1})$  und  $v = (w_1, ..., w_n, s)$ , so ist

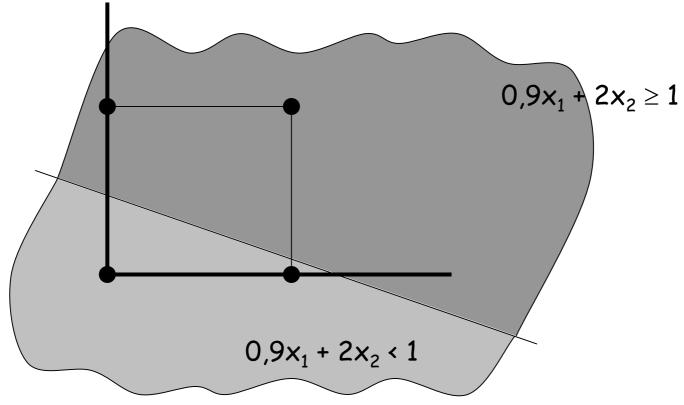
o = 0, falls 
$$yv = x_1 w_1 + ... + x_n w_n - 1s < 0$$

o = 1, sonst

# Geometrische Interpretation

- Gleichung xw = s beschreibt eine Hyperebene im n -dimensionalen Raum
- · Beispiel:

$$0.9x_1 + 2x_2 = 1$$

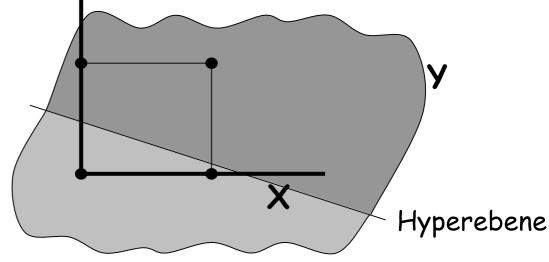


#### Lernen der Gewichte

- · Gegeben ist eine Menge von Beispielen
- Überwachte Lernaufgabe: Daten sind disjunkt in zwei Mengen X,Y geteilt
- gesucht ist der Gewichtsvektor  $w(w_1, ..., w_n)$ , so daß eine Hyperebene spezifiziert wird, die beide Mengen X und Y voneinander trennt

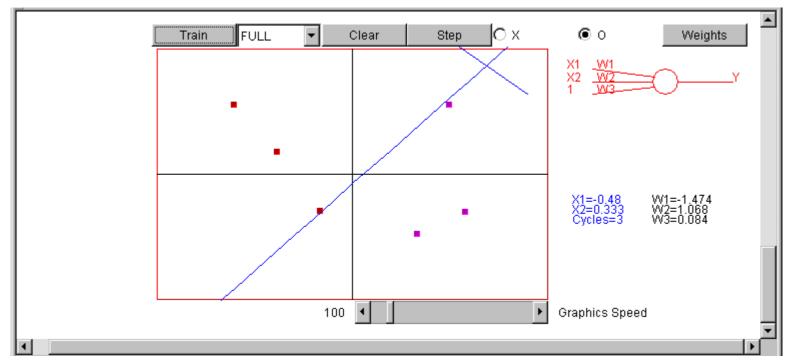
Mengen müssen linear trennbar sein, damit die Aufgabe

lösbar ist.



# Delta-Regel

Menge	ΣX	Menge Y		
$x_1$	<b>X</b> <sub>2</sub>	$\mathbf{x}_1$	<b>X</b> <sub>2</sub>	
0,552	-0,605	0,552 -0,304	0,497	
0,176	-0,384	-0,304	0,579	
-0,296	-0,164	-0,48	0,333	



# Delta-Regel

- beim Training werden die Beispiele dem Netz als Input präsentiert
- Output ist für die Beispiele bekannt
  --> überwachte Lernaufgabe (supervised)

(hier: liegt Beispiel in X oder Y?)

• Soll und Ist-Output werden verglichen und bei Diskrepanz werden Schwellenwert und Gewichte nach folgender Delta-Regel angepasst:  $(w_0 = -s, x_0 = 1)$ 

η Lernrate

$$w_{i,neu} = w_{i,alt} + \eta x_i * (Output_{soll} - Output_{ist})$$

# Delta-Regel

Algorithmus mit Lernrate  $\eta$  = 1, als Output nur 0 und 1 möglich!

Start: Der Gewichtsvektor  $\mathbf{w}_0$  wird zufällig generiert.

Setze t:=0.

Testen: Ein Punkt x in  $X \cup Y$  wird zufällig gewählt

Falls  $x \in X$  und  $w_t \cdot x > 0$  gehe zu Testen

Falls  $x \in X$  und  $w_{+} \cdot x \leq 0$  gehe zu Addieren

Falls  $x \in Y$  und  $w_t \cdot x < 0$  gehe zu Testen

Falls  $x \in Y$  und  $w_t \cdot x \ge 0$  gehe zu Subtrahieren

Addieren: Setze  $w_{t+1} = w_t + x$ .

Setze t := t + 1. Gehe zu Testen

Subtrahieren: Setze  $w_{t+1} = w_t - x$ .

Setze t := t + 1. Gehe zu Testen

# Delta-Regel - Beispiel

t x X/Y Vhs. Error zu\_addieren/subtr.

neue\_Gewichte

	i	t	$a_v$	e	$\Delta W(u_1,v)$	$\Delta W(u_2,v)$	$\Delta \theta$	$W(u_1,v)$	$W(u_2,v)$	θ
	0.0	0	0	0	0	0	0.	0	0	0
1 Ek-	0 1	0	0	0	0	0	0	0	0	0
1. Epoche	10	0	0	0	0	0	0	0	0	0
	11	1	0	1	1	1	-1	1	1	-1
	0 0	0	1	-1	0	0	1	1	1	0
2. Epoche	0 1	0	1	-1	0	-1	1	1	0	1
2. Epoche	10	0	0	0	0	0	0	1	0	1
	1 1	$1_{1}$	0	1	1	1	-1	2	1, .	0
3	0 0	0	0	0	0	0	0	2	1	0
3. Epoche	0 1	0	1	-1	0	-1	1	2	0	1
3. Epoche	10	0	1	1-1	-1	0	1	1	0	2
engerner to	1.1	1	0	1	1	1	-1	2	1	1
1 1 1 1 1	0.0	0	0	0	0	0	0	2	1	1
4. Epoche	0 1	0	0	0	0	0	0	2	1	1
4. Epoche	10	0	1	-1	-1	0	1	1	1	2
Total variables	11	1	0	1	1	1	-1	2	2	1
	0 0	0	0	0	0	0	0	2	2	1
5. Epoche	0.1	0	1	-1	0	-1	- 1	2	1	2
o. Epoche	10	0	0	0	0	0	0	2	1	2
	1 1	1	1	0	0	0	0	2	1	2
6. Epoche	0 0	0	0	0	0	0	0	2	1	2
	0 1	0	0	0	0	. 0	0	2	1	2
o. Epoche	10	0	0	0	0	0	0	2	1	2
	1 1	1	1	0	0	0	0	2	1	2

entnommen Nauk, Kruse, S. 50

# Konvergenz und Korrektheit der Delta-Regel

Satz: Wenn das Perzeptron eine Klasseneinteilung überhaupt lernen kann, dann lernt es diese mit der

Delta-Regel in endlich vielen Schritten.

Problem: Falls das Perzeptron nicht lernt, kann nicht unterschieden werden, ob nur noch nicht genügend Schritte vollzogen wurden oder ob das Problem nicht lernbar ist (keine obere Schranke)

# Konvergenz und Korrektheit der Delta-Regel

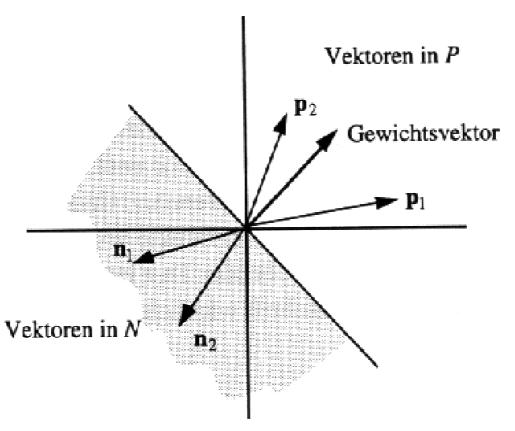


Abb. 4.8 Positiver und negativer Halbraum im Eingaberaum

entnommen Rojas, S. 84

# Konvergenz und Korrektheit der Delta-Regel

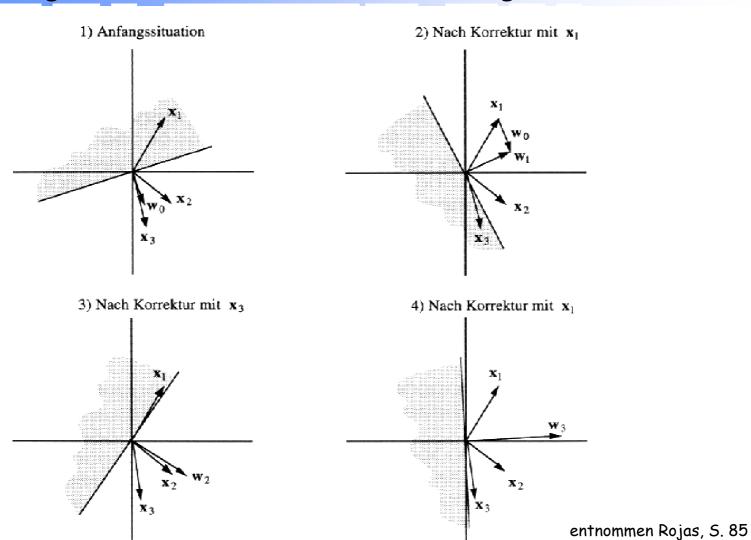
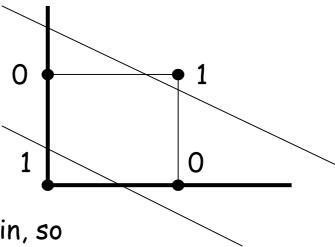


Abb. 4.9 Konvergenzverhalten des Lernalgorithmus

#### XOR-Problem

Logisches AND ist linear separierbar

Logisches XOR ist nicht linear separierbar



Führt man weitere Hyperebenen ein, so kann man auch hier die Klassen unterscheiden. ==> Realisierung durch Zwischenschichten

- Perzeptron kann z.B. verschiedene <u>Boolsche Funktionen</u> repräsentieren: AND, OR, NAND, NOR (aber: XOR ist <u>nicht</u> darstellbar)
  - Realisierung von AND:
    - · wähle Perzeptron mit zwei Eingabegrößen
    - wähle  $w_0 = -0.8$ ,  $w_1 = w_2 = 0.5$
  - AND, OR sind Spezialfall sogenannter "m-aus-n Funktionen" (m-of-n functions): wenigstens m der n Eingabegrößen müssen 1 sein, hier:
    - AND: m=1 (wenigstens 1 Eingabegröße ist 1)
    - OR: m=n (alle Eingabegrößen sind 1)

# Agenda

- 1. Einführung
- 2. Einfaches Perzeptron
- 3. Multi Layer Perzeptron
  - Vektorschreibweise der Deltaregel
  - Schichten des MLP
  - Backpropagation
  - Probleme der Backpropagation
  - Varianten der Backpropagation
- 4. Beispiel
- 5. Rekurrente Netze
- 6. Entwurfsüberlegungen

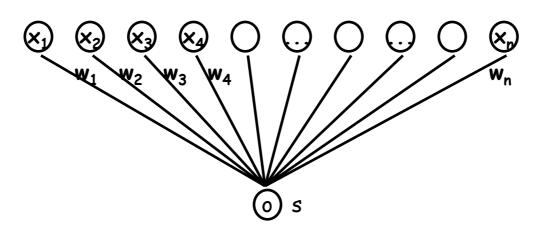
# Erinnerung Perzeptron

$$\Delta w_i := \eta \cdot x_i \cdot e$$
 (Synapsenveränderung)

Schwellwert:

• Netz-Output:  $o = \theta(x_1w_1 + x_2w_2 + x_3w_3 + ... + x_nw_n - s)$ • erwarteter Output: t (target) • gemachter Fehler: e = t - o (error)

· Lernkonstante:



# Vektorschreibweise der Delta Regel

• Aktivität der Input-Schicht:

$$\underline{\mathbf{x}} = (\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \dots, \mathbf{x}_n, \mathbf{1})$$

Gewichtsvektor (einschl. Schwellwert):

$$\underline{\mathbf{w}} = (\mathbf{w}_1, \mathbf{w}_2, \mathbf{w}_3, \dots, \mathbf{w}_n, -\mathbf{s})$$

Aktivität des Ausgabeneurons:

$$o = \theta(\underline{\mathbf{x}} \cdot \underline{\mathbf{w}})$$

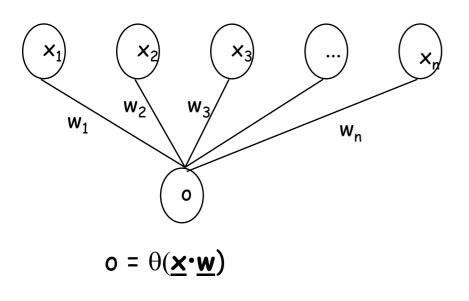
• Fehler des Ausgabeneurons (t = erwarteter Wert):

$$e = t-o$$

Gewichtsänderung:

$$\Delta \mathbf{w} := \eta \cdot \mathbf{e} \cdot \mathbf{x}$$

# Delta Regel als Ableitung für Perzeptron



### Fehlergradient:

$$F = (o - t)^{2} = (\theta(\underline{\mathbf{x}} \cdot \underline{\mathbf{w}}) - t)^{2}$$

$$\partial F / \partial w_{i} = \partial(o - t) / \partial w_{i}$$

$$= \partial(\theta(\underline{\mathbf{x}} \cdot \underline{\mathbf{w}}) - t)^{2} / \partial w_{i}$$

$$= 2 \theta'(\underline{\mathbf{x}} \cdot \underline{\mathbf{w}}) (o - t) x_{i}$$

Die Delta Regel kann als Gradientenabstieg mit (variablen) Lernfaktor interpretiert werden:

$$\Delta w_i = \eta \text{ (o-t) } x_i \text{ mit } \eta = 2 \theta' (\underline{\mathbf{x}} \cdot \underline{\mathbf{w}})$$

(unter derAnnahme:  $\theta$  ist diff.-bar)

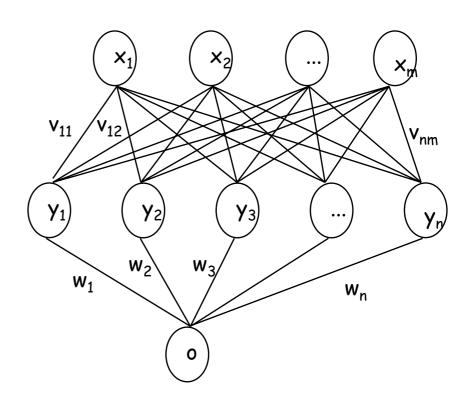
# 2 Layer Perzeptron

Input-Vektor <u>x</u>
Gewichtsmatrix <u>v</u>

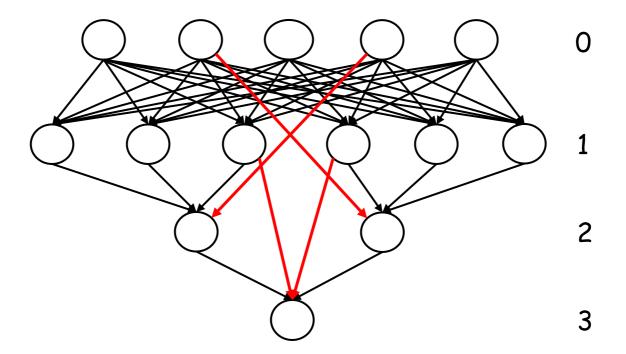
Aktivitätsvektor <u>y</u> Gewichtsvektor <u>w</u>

Output o  $\mathbf{y} = \theta(\mathbf{v} \cdot \mathbf{x})$ 

$$o = \theta(\underline{\mathbf{w}} \cdot \underline{\mathbf{y}})$$



 Netze ohne Rückkopplungen können stets in die Form eines Schichtennetzwerkes gebracht werden, wobei aber durchaus Verbindungen eine oder mehrere Schichten überspringen dürfen

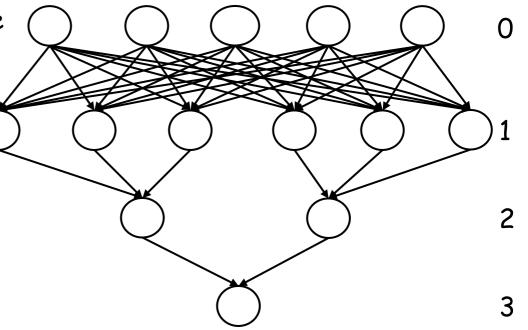


 Wenn in einem Schichtennetz keine Verbindungen existieren, die eine oder mehrere Schichten überspringen, spricht man von einem strengen Schichtennetz.

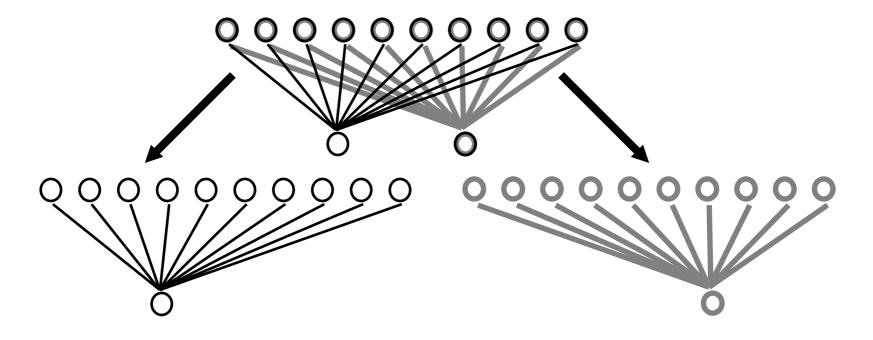
 In diesem Fall ist jede Schicht durch die Länge des Weges zu ihren Neuronen gekennzeichnet.

 Die input-Schicht hat Weglänge 0.

 Die Output-Schicht hat maximale Weglänge.

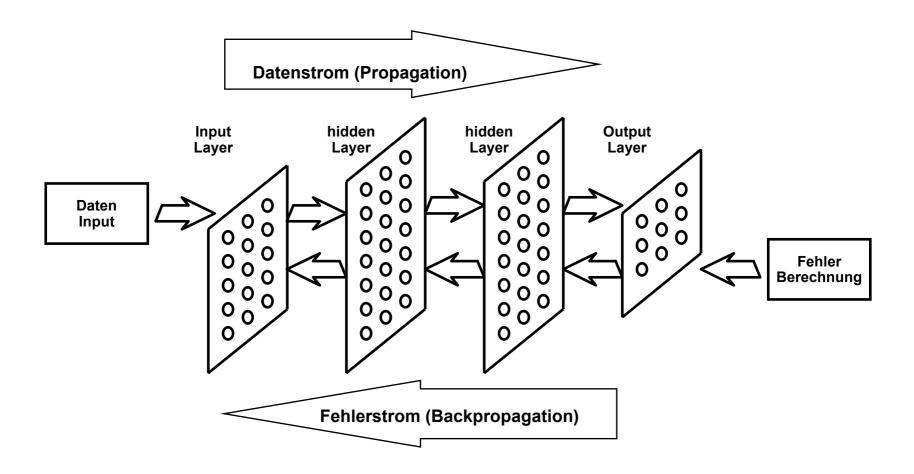


Die beiden Neuronen der 2. Schicht beeinflussen einander nicht, deshalb können sie voneinander getrennt betrachtet werden.



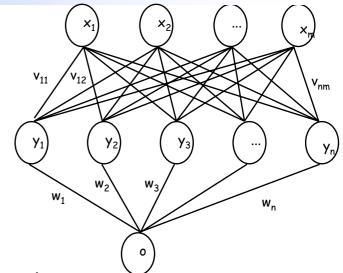
Bei mehrschichtigen Netzen geht die Unabhängigkeit verloren, d.h. sie können nicht so getrennt werden.

# Backpropagation



Fehlerfunktion F (mittlerer quatratischer Fehler) für das Lernen:

$$= \frac{1}{2} \sum_{d \in D} \left( t_d - o_d \right)^2$$



wobei gilt:

D Menge der Trainingsbeispiele

 $t_d$  korrekter Output für  $d \in D$ 

 $o_d^a$  berechneter Output für  $d \in D$ 

Anders geschrieben ist

$$F = (o - t)^2 = (\theta(\underline{\mathbf{w}} \cdot \underline{\mathbf{y}}) - t)^2$$

Die Gewichte müssen so angepasst werden, daß der Fehler minimiert wird. Dazu bietet sich das Gradientenabstiegsverfahren an.

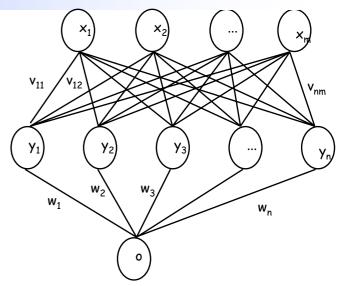
# Fehlergradient für wi lautet:

$$\partial F/\partial w_{i} = \partial(o-t)^{2}/\partial w_{i}$$

$$= \partial(\theta(\underline{w}\cdot\underline{y}) - t)^{2}/\partial w_{i}$$

$$= 2\cdot(o-t)\cdot\theta'(\underline{w}\cdot\underline{y})\cdot\partial(\underline{w}\cdot\underline{y})/\partial w_{i}$$

$$= 2\cdot(o-t)\cdot\theta'(\underline{w}\cdot\underline{y})y_{i}$$



# Fehlergradient für v<sub>ii</sub> lautet:

$$\partial F/\partial v_{ij} = \partial F/\partial y_{i} \cdot \partial y_{i}/\partial v_{ij}$$

$$= \partial F/\partial y_{i} \cdot \partial \theta(\underline{v}_{i} \cdot \underline{x})/\partial v_{ij}$$

$$(Fehler von Neuron i)$$

$$= \partial F/\partial y_{i} \cdot \theta'(\underline{v}_{i} \cdot \underline{x}) \cdot x_{j}$$

$$= \partial F/\partial o \cdot \partial o/\partial y_{i} \cdot \theta'(\underline{v}_{i} \cdot \underline{x}) \cdot x_{j}$$

$$= \partial F/\partial o \cdot \partial \theta(\underline{w} \cdot \underline{y})/\partial y_{i} \cdot \theta'(\underline{v}_{i} \cdot \underline{x}) \cdot x_{j}$$

$$= \partial F/\partial o \cdot \theta'(\underline{w} \cdot \underline{y}) \cdot w_{i} \cdot \theta'(\underline{v}_{i} \cdot \underline{x}) \cdot x_{j}$$

$$= \partial F/\partial o \cdot \theta'(\underline{w} \cdot \underline{y}) \cdot w_{i} \cdot \theta'(\underline{v}_{i} \cdot \underline{x}) \cdot x_{j}$$

$$= 2 \cdot (o-t) \cdot \theta'(\underline{w} \cdot \underline{y}) \cdot w_{i} \cdot \theta'(\underline{v}_{i} \cdot \underline{x}) \cdot x_{j}$$
Input

Fehler bei Info von

der Ausgabe Zwischenschicht

Info von Inputschicht

- Die schichtweise Berechnung der Gradienten ist auch für mehr als zwei Schichten möglich
- Fehler wird schichtenweise nach oben propagiert (Back-Propagation)
- · allgemein gilt für den Output-Fehler e eines Neurons j  $e_j = \partial F/\partial x_j$
- Gewichtsänderung  $\Delta w_{ik} = a \cdot e_k \cdot \theta'(a_k) \cdot x_i$

### Backpropagation Algorithmus

- Wähle ein Muster x aus der Menge der Trainingsbeispiele D aus
- präsentiere das Muster dem Netz und propagiere es
- anhand der resultierenden Neuronenaktivität wird der Fehler Fermittelt
- die Fehlerinformationen werden durch das Netz zurückpropagiert (Backpropagation)
- die Gewichte zwischen den einzelnen Schichten werden so angepasst ( $\Delta w_{ik}$ ), dass der mittlere Ausgabefehler für das Muster sinkt

# Backpropagation Algorithmus

 Betrachtet man den Fehler des Netzes als Funktion aller Gewichte w, so kann man zeigen, daß bei jedem Schritt der Fehler kleiner wird. Dies erfolgt unabhängig vom gewählten Netz, also unabhängig von w.

• 
$$e(w) = (o - t)^2 = (t - \theta (w \cdot x))^2$$

· es gilt bei jedem Schritt:

$$e(\mathbf{w} + \Delta \mathbf{w})^2 < e^2$$

# Backpropagation Algorithmus

- · Die Gewichtsänderungen können auf zwei Arten erfolgen:
  - Online Training: jedes Gewicht wird sofort angepaßt (folgt nur im Mittel dem Gradienten)
  - Batch Verfahren: es werden alle Datensätze präsentiert, die Gewichtsänderung des Gewichtes berechnet, summiert und dann erst angepaßt (entspricht dem Gradienten über dem Datensatz)

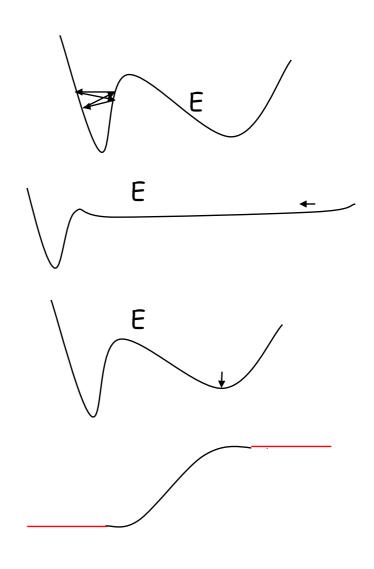
# Probleme des Backpropagation Algorithmus

 Oszillation in engen Schluchten

 Stagnation auf flachen Plateaus

· lokale Minima

• Flat Spots ( $\theta'(a_j) \approx 0$ ) kaum Veränderung im Training



# Varianten des Backpropagation Algorithmus

- · Änderung der Fehlerfunktion
- · Momentum-Term einführen
- Weight Decay
- Lernkonstanten und Output Funktion können für jede
   Schicht oder für jedes Neuron einzeln festgelegt werden
- · Lernkonstanten können dynamisch angepaßt werden
- · Gewichtsanpassung modifizieren (Manhatten Training)
- Flat Spots eleminieren
- Lernrate  $\eta$  wird an die Gewichte und den Gradienten gekoppelt
- Konvergenz durch Nutzung der zweiten Ableitung beschleunigen

#### – Bemerkungen:

- jede <u>Boolesche Funktion</u> kann durch derartiges Netz repräsentiert werden
- <u>bel. Funktionen</u> können durch ein ANN mit drei Schichten bel. genau <u>approximiert</u> werden
- Hypothesenraum ist kontinuierlich im Gegensatz z.B. zum diskreten Hypothesenraum von Entscheidungsbaumverfahren
- ANN kann für interne Schicht sinnvolle Repräsentationen lernen, die im vorhinein <u>nicht</u> bekannt sind; dies ist Gegensatz zu z.B. ILP - Verfahren mit <u>vorgegebenem</u> Hintergrundwissen

### VI.3.3.3 Abbruchbedingungen und Überspezialisierung

- Beschreibung des BACKPROPAGATION Algorithmus läßt Abbruchbedingung offen
  - Algorithmus solange zu iterieren, bis Fehler für die Trainingsbeispiele unter einem vorgegebenen Schwellwert liegt, ist <u>schlechte</u> Strategie, da Algorithmus zur <u>Überspezialisierung</u> (overfitting) neigt ⇒ ( vgl. Entscheidungsbäume)
  - verwende <u>separate</u> Menge von Beispielen zur <u>Validierung</u> (validation set); iteriere solange, bis Fehler für Validierungsmenge <u>minimal</u> ist
  - <u>aber</u>: Fehler auf Validierungsmenge muß nicht monoton fallen! (im Gegensatz zu Fehler auf Trainingsmenge. Vgl. auch Abbildung 8 a/b)

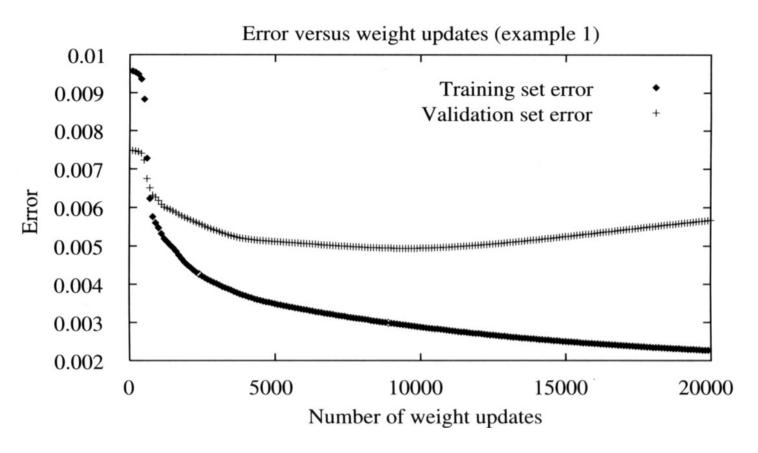


Figure 8a: Plots of error E as a function of the number of weight updates, for two different robot perception tasks (Mitchell 1997)

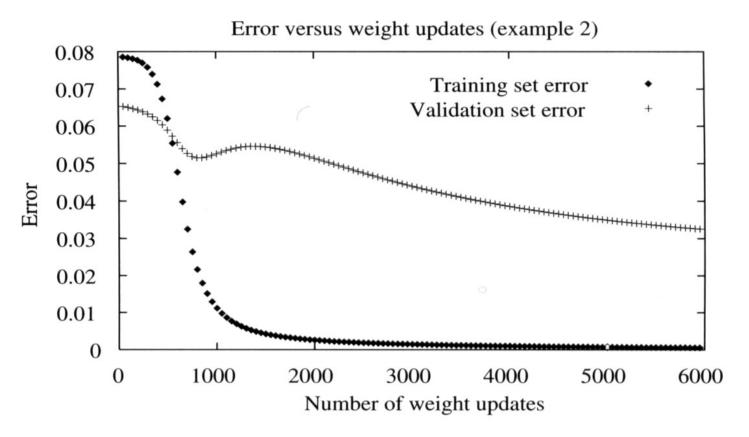


Figure 8b: Plots of error E as a function of the number of weight updates, for two different robot perception tasks (Mitchell 1997)

- Alternativ kann <u>k-fache Kreuzvalidierung</u> (<u>k-fold cross-validation</u>) verwendet werden:
  - unterteile Menge der Trainingsbeispiele in k gleich große disjunkte Teilmengen
  - verwende der Reihe nach jeweils eine andere Teilmenge als Validierungsmenge und die restlichen (k - 1) Teilmengen als Trainingsmenge
  - für jede Validierungsmenge wird "optimale" Anzahl i von Iterationen bestimmt (d.h. mit kleinstem Fehler für Beispiele aus Validierungsmenge)
  - <u>Mittelwert</u> von i über alle k Trainingsphasen wird letztendlich verwendet, um geg. ANN mit <u>allen</u> Trainingsbeispielen zu trainieren

# Zusammenfassung Backpropagation Algorithmus

- Backpropagation stellt trotz der beschriebenen Probleme einen bedeutenden Fortschritt dar.
- Aufgaben, wie das "exklusive oder", die mit einem Perzeptron nicht gelöst werden können, sind nun lösbar.
- Die Lösung wird durch die innere Schicht ermöglicht, die eine Umkodierung der Eingabevektoren vornimmt.

# Agenda

- 1. Einführung
- 2. Einfaches Perzeptron
- 3. Multi Layer Perzeptron
- 4. Beispiel
- 5. Rekurrente Netze
- 6. Entwurfsüberlegungen

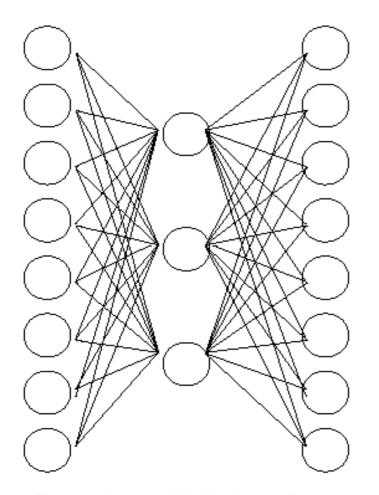
### Voraussetzungen

- Netz aus drei Schichten
- · Ein- und Ausgabe enthalten jeweils N Neuronen
- · die mittlere Schicht enthält nur M < N Neuronen

### Lernaufgabe:

- Wenn Neuron n in der Eingabeschicht aktiviert wird, dann soll auch Neuron n auf der Ausgabeseite aktiviert werden.
- Die Aufgabe wäre trivial, wenn die mittlere Schicht ebenfalls N Neuronen hätte.

==> mittlere Schicht stellt Flaschenhals dar



Input	H	Hidden		Output
		Values		
10000000	> .89	.04 .08	>	10000000
01000000	> .15	.99 .99	>	01000000
00100000	> .01	.97 .27	>	00100000
00010000	> .99	.97 .71	>	00010000
00001000	> .03	.05 .02	>	00001000
00000100	> .01	.11 .88	>	00000100
00000010	> .80	.01 .98	>	00000010
00000001	> .60	.94 .01	>	0 0 0 0 0 0 0 1

Figure 7: Learned Hidden Layer Representation (Mitchell 1997)

### Lösung:

 Das Netz muss eine geeignete Kodierung für die Representation der aktivierten Neuronen in der mittleren Schicht finden.

### Vorgehen:

- · Gewichte werden zufällig gewählt
- 0.9 wird als 1 und 0.1 als 0 interpretiert (Konvergenzproblem)

### Ergebnis:

· 3 bit Binärcode für die innere Repräsentation

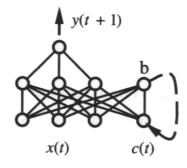
	000	
_ · · · · · · ·		
· · ·		
	• • • • • • • • • • • • • • • • • • •	
• • •	·	
		entnommen Ritter, S. 59

Frage: Welche interne Datenrepräsentationen sind erforderlich, um eine gegebene Aufgabe durch ein massiv paralleles Netzwerk lösen zu können?

- Backpropagation erlaubt das Trainieren von Netzen auf unterschiedliche Probleme.
- Die Analyse der "Neuroanatomie" gibt Einblicke in die Netzstruktur und ermöglicht ein besseres Verständnis auch von biologischen Netzwerken.

#### VI.3.3.5 Rekurrente Netzwerke

- Verarbeitung von <u>Zeitreihen</u> erfordert die Berechnung von Ausgabewerte zum Zeitpunkt (t+1) in Abhängigkeit der Ausgabewerte t, t-1, ...
  - z.B. bei der Prognose von Aktienkursen
- Zeitfenster in die Vergangenheit wird benötigt:



(b) Recurrent network

 Eingabewert c(t) ist gleich dem Wert der internen Einheit b zum Zeitpunkt (t-1)

- interne Einheit b sammelt Historie über Eingabewerte für das Netzwerk auf:
  - Wert von b(t) ist abhängig von aktuellen Eingabewerten x(t) und c(t), d.h. b(t-1)
  - damit ist Wert von b(t) indirekt abhängig von Eingabewerten x(t), x(t-1), ...
- im Prinzip kann ein rekurrentes Netzwerk mit einer leichten Variation des Backpropagation-Algorithmus trainiert werden (siehe z.B. Mozer 1995)

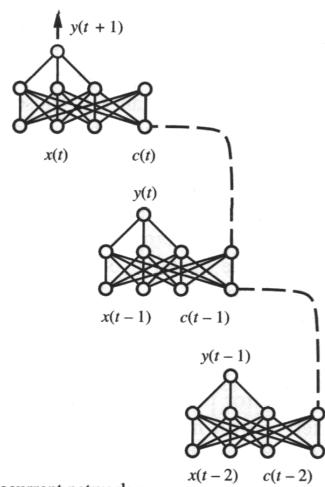


Figure 9: Recurrent Networks (Mitchell 1997)

(c) Recurrent network unfolded in time

#### VI.3.3.4 Entwurfsüberlegungen für ANNs

#### – Eingabecodierung:

 welche Repräsentation der vorliegenden Trainingsbeispiele ist günstig - Repräsentation muß <u>fixe</u> Anzahl von Eingabeelementen ermöglichen

#### – <u>Ausgabecodierung</u>:

- welche Repräsentation ist günstig
- sofern n diskrete Ausgabewerte berechnet werden sollen, bietet sich <u>1-aus-n-Codierung</u> an (<u>1-of-n output encoding</u>): jeder Ausgabewert wird durch 1 Ausgabeeinheit repräsentiert
  - viele Kanten ermöglichen eine <u>flexible</u> Repräsentation der Ausgabefunktion über ihre zugehörigen Gewichte
  - <u>Differenz</u> zwischen dem höchsten und dem zweithöchsten Ausgabewert kann als Maß für die <u>Güte</u> der berechneten Ausgaben verwendet werden

- Struktur der internen Schicht
  - wieviele Elemente soll man für die interne Schicht wählen?
  - "gewisse minimale" Anzahl von Elementen wird benötigt, um Zielfunktion genügend genau zu lernen
  - Kreuzvalidierung kann Überspezialisierung durch zu große Anzahl interner Elemente vermeiden

<u>aber</u>: bisher fehlt klare Methodik für systematischen Entwurf von ANNs (siehe z.B. (Anders 1997))

#### Weiteres Beispiel: Mustererkennung mit Hopfield-Netzen

→ Demo Hopfield.html

Bei Hopfield-Netzen sind die Ausgabe-Neuronen wieder mit den Eingabe-Neuronen verbunden.

Die Eingabe ist nun durch eine initiale Eingabe kodiert.

Das Netz versucht, durch Backpropagation in einen stabilen Zustand zu kommen.

In der Demo ist jedes Feld der Matrix einem von 64 Ausgabeneuronen zugeordnet.