Vorlesung Künstliche Intelligenz Wintersemester 2008/09

Teil III: Wissensrepräsentation und Inferenz

Kap.5: Neuronale Netze

- Dieses Kapitel basiert auf Material von Andreas Hotho
- Mehr Details sind in der Vorlesung "Neuronale Netze" von Prof. Werner zu finden.



- 1. Einführung & Grundbegriffe
 - Motivation & Definition
 - Vorbild Biologie
- 2. Einfaches Perzeptron
- 3. Multi-Layer-Perzeptron
- 4. Vor- und Nachteile Neuronaler Netze



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.

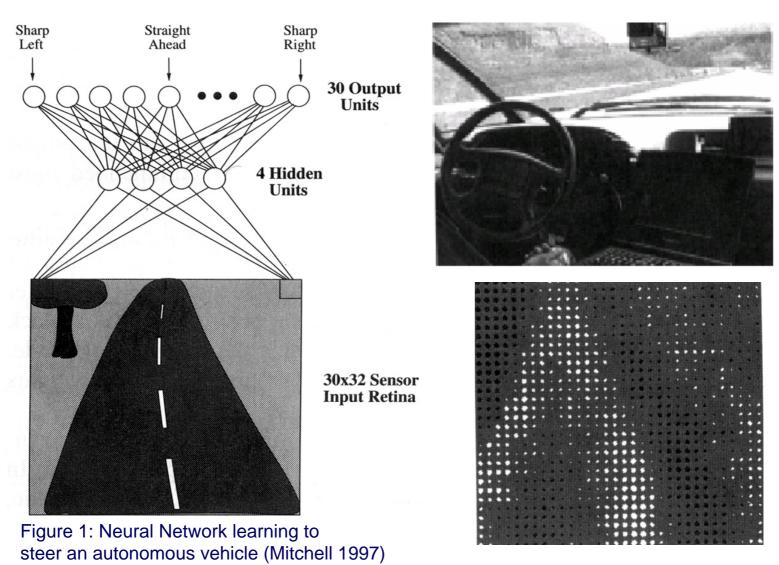
(Kohonen 84)



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







<u>Charakteristische</u> Eigenschaften von Problemen, die mit BACKPROPAGATION-ANNs gelöst werden können:

- <u>Trainingsbeispiele</u> sind repräsentiert durch <u>Attribut-Wert-</u> Paare
 - · Werte können reellwertig sein
- + zu generierende Funktion (target function) kann ...
 - · <u>diskrete</u> Funktionswerte
 - · reellwertige Funktionswerte oder
 - · <u>Vektor</u> solcher Funktionswerte

haben

- + Trainingsbeispiele dürfen fehlerhaft sein
- lange 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

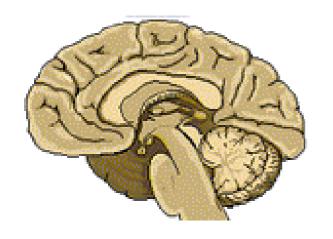
Es werden

- Aufbauphase (Topologie),
- · Trainingsphase (Lernen) und
- · Arbeitsphase (Propagation) unterschieden.

(Die Phasen können auch überlappen.)

Vorbild Biologie

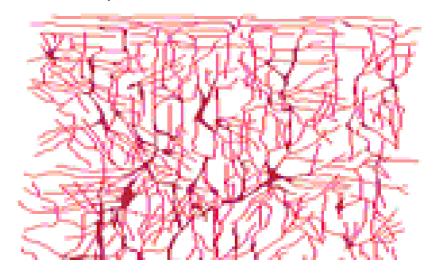
Gehirn



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

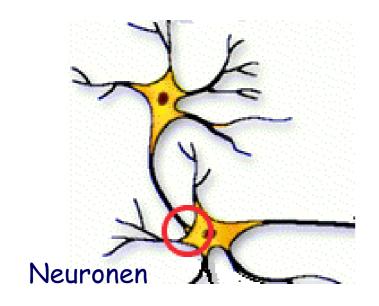
Geflecht aus Neuronen

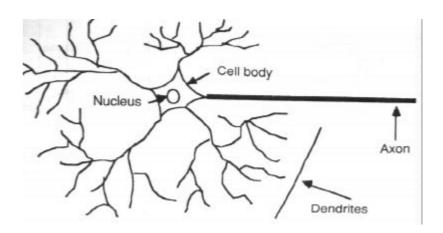
Das Gehirn besteht aus ca. 10¹¹ Neuronen, die jeweils mit max. 10⁴ anderen Neuronen durch insgesamt ca. 10¹³ Synapsen verschaltet sind.

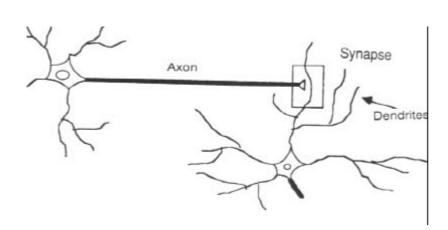


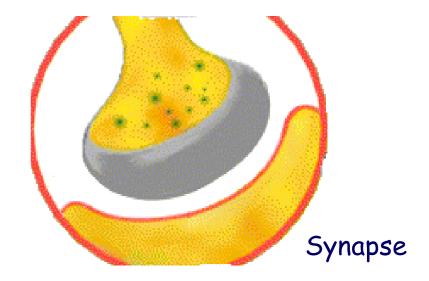
Vorbild Biologie







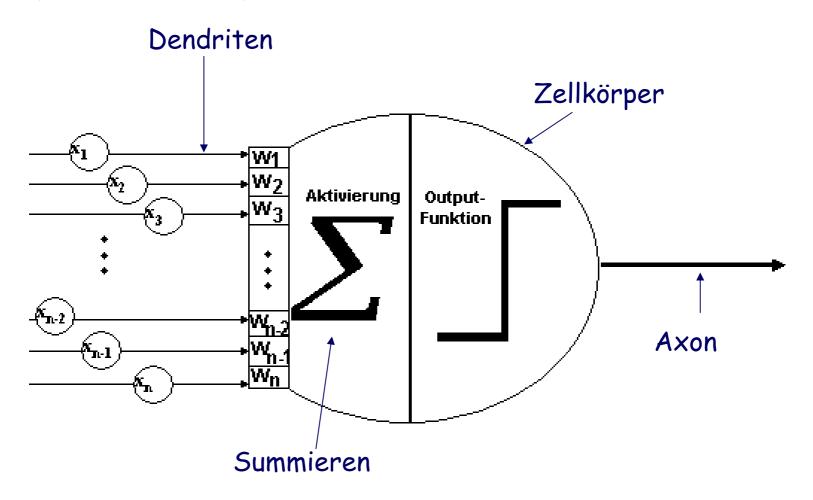




Vorbild Biologie

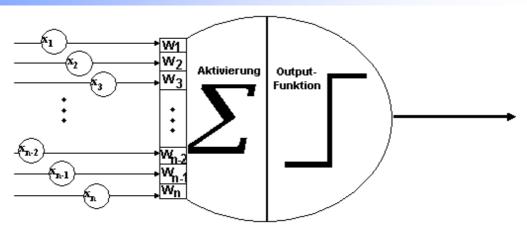


Vom natürlichen zum künstlichen Neuronalen Netz



Aktivierungsfunktionen





Es gibt mehrere 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: $a_1 = x_1^* w_1$, $a_2 = x_2^* w_2, ..., a_n = x_n^* w_n$.

Aktivierungsfunktionen



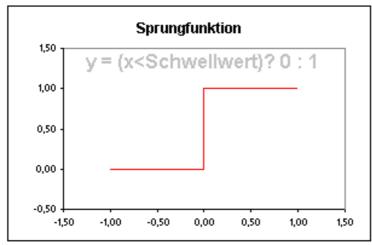
- Allen gemeinsam ist die **Gewichtung** der Inputs. Die Inputs $x_1, x_2, ..., x_n$ werden stets mit ihren (Synapsen-)Gewichten $w_1, w_2, ..., w_n$ multipliziert: $a_1 = x_1^* w_1$, $a_2 = x_2^* w_2, ..., a_n = x_n^* w_n$.
- 1. Die am häufigsten angewandte Regel ist die **Skalarprodukt-Regel**: Die gewichteten Inputs $a_1,a_2,...,a_n$ werden zur Aktivität des Neurons aufaddiert:

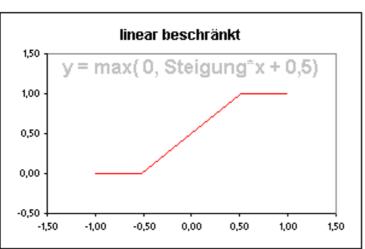
$$a = a_1 + a_2 + ... + a_n$$

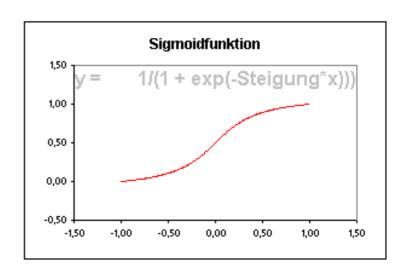
2. Sehr häufig ist ebenfalls die Winner-take-all-Regel: bei der die Aktivität a zunächst nach der Skalarproduktregel ermittelt wird, dann aber mit allen Aktivitäten in derselben Schicht verglichen wird und auf 0 herabgesetzt wird, wenn ein anderes Neuron höhere Aktivität hat.

Outputfunktionen









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

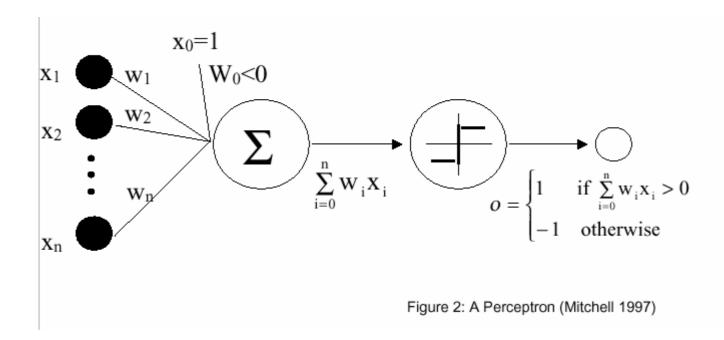
Agenda



- 1. Einführung
- 2. Einfaches Perzeptron
 - Definition Perzeptron
 - Geometrische Interpretation
 - Lernen: Delta-Regel
 - XOR-Problem
- 3. Multi-Layer-Perzeptron
- 4. Vor- und Nachteile Neuronaler Netze

Perzeptron





Das einfache Perzeptron besteht aus einem einzelnen Neuron. • 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
$$yv = x_1 w_1 + ... + x_n w_n - 1s < 0$$

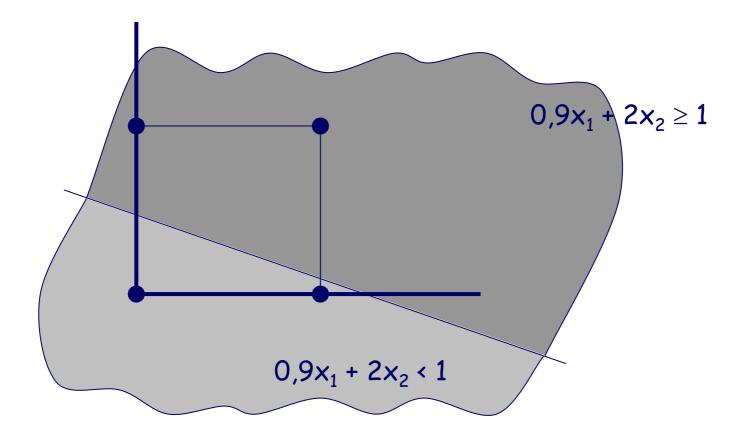
o = 1, sonst

Geometrische Interpretation



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

$$0.9x_1 + 2x_2 = 1$$



Lernen der Gewichte

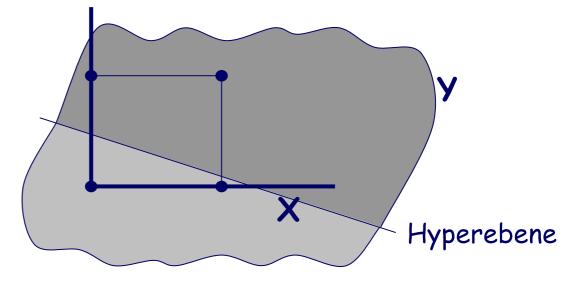


- · Gegeben ist eine Menge von Beispielen.
- Überwachte Lernaufgabe: Daten sind disjunkt in zwei Mengen X,Y geteilt
- Gesucht ist der Gewichtsvektor $w(w_1, ..., w_n)$, so dass eine Hyperebene spezifiziert wird, die die Mengen X und Y voneinander trennt.

Mengen m

üssen linear trennbar sein, damit die Aufgabe

lösbar ist.



Konvergenz und Korrektheit der Delta-Regel



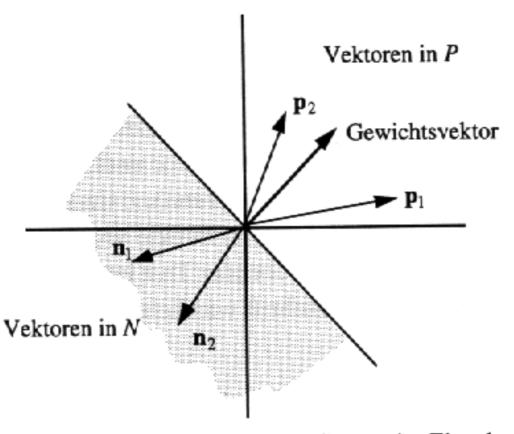


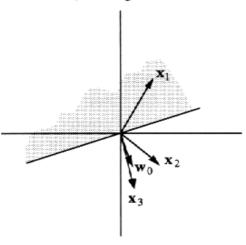
Abb. 4.8 Positiver und negativer Halbraum im Eingaberaum

entnommen Rojas, S. 84

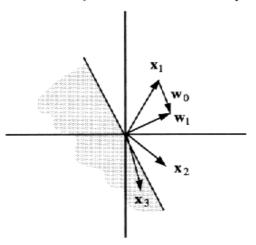
Konvergenz und Korrektheit der Delta-Regel



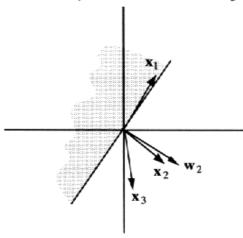
1) Anfangssituation



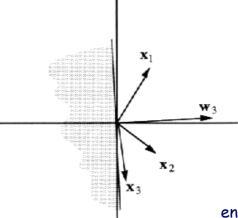
2) Nach Korrektur mit x₁



3) Nach Korrektur mit x3



4) Nach Korrektur mit x₁



entnommen Rojas, S. 85

Abb. 4.9 Konvergenzverhalten des Lernalgorithmus

Delta-Regel

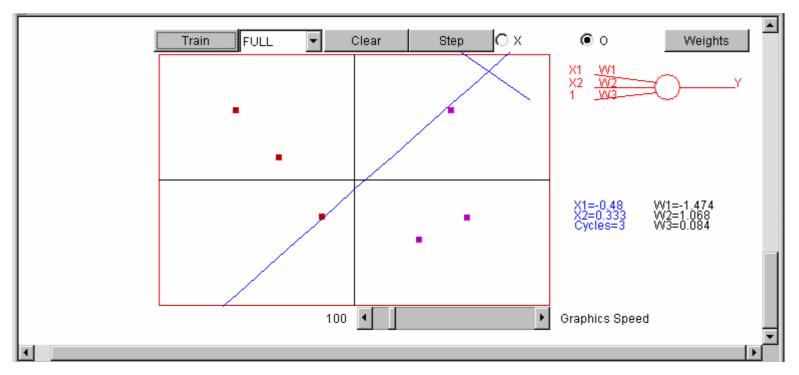


Menge	X
-------	---

$\underline{\hspace{1cm}}^{\hspace{1cm}} \hspace{1cm} \hspace{1cm}\hspace{1cm} \hspace{1cm} 1$	X ₂
0,552	-0,605
	-0,384
-0,296	-0,164

Menge Y

X ₁	X ₂
0,552	0,497
-0,304	0,579
-0,48	0,333



 \rightarrow Demo (http://neuron.eng.wayne.edu/java/Perceptron/New38.html)

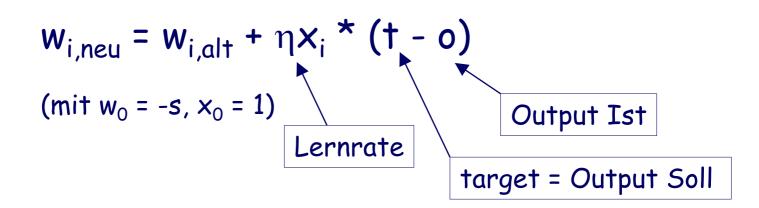
Delta-Regel



- 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.
 Bei Diskrepanz werden Schwellenwert und Gewichte nach folgender Delta-Regel angepasst:



Delta-Regel



Annahme hier: - Algorithmus mit Lernrate $\eta = 1$

- als Output nur 0 und 1 möglich (d.h. Trennung

von zwei Klassen wird gelernt)

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

Setze j:=0. // j=Zeit

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

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

Falls $x \in X$ and $w_j \cdot x \le 0$ gehe zu Addieren

Falls $x \in Y$ und $w_j \cdot x \le 0$ gehe zu Testen

Falls $x \in Y$ und $w_i \cdot x > 0$ gehe zu Subtrahieren

Addieren:

Setze $w_{j+1} = w_j + x$. Setze j := j + 1. Gehe zu Testen

Subtrahieren: Setze $w_{j+1} = w_j - x$. Setze j := j + 1. Gehe zu Testen

Beispiel für Anwendung der Delta-Regel



Wir wollen das logische Und lernen.

i	t
0.0	0
0 1	0
10	0
11	1

Start:

Der Gewichtsvektor w_0 wird "zufällig" generiert: w_1 :=0, w_2 :=0, Schwellwert $\theta = -w_0$:=0

Zeit $|x_1 x_2| t$ | 0 | Error | zu_addieren/subtr.

neue_Gewichte

		_	and the same of the same							
	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	Q	0
1. Epoche	0 1	0	0	0	0	0	0	0	0	0
т. Ероспе	1 0	0	0	0	0	0	0	0	\0 \	0
	1 1	1	0	1	1	1	-1	1	ì	\mathbf{M}
	0 0	0	1	-1	0	0	1	1	1	0
9 Enoche	0 1	0	1	-1	0	-1	1	1	0	1

In unserer Notation:

$$w_1 \ w_2 - w_0 = s$$

Beispiel für Anwendung der Delta-Regel

Wir wollen das logische Und lernen.

	i	t	Γ
1	0 0	0	Γ
	0 1	0	l
	10	0	
	1 1	1	

Falls $x \in X$ und $w_t \cdot x \leq 0$ gehe zu Addieren

Zeit	$ x_1 x_2 $	t	o Error zu_addieren/subtr.					neue_	Gewichte	- 1
	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. Epoche	0 1	0	0	0	0	0	0	0	0	0
1. Epoche	1 0	0	0	0	0	0	0	0	0	0
	1 1	1	0	1	1	1	-1		1	-1
	0 0	0	1	+1	0	0	1	1	1	0
2 Enoche	0 1	Ø	X	-1	0	-1	1	1	0	1
		<	1x0	+ 1x0	$-1_{const} \times 0 = 0$					

Zeit

Beispiel für Anwendung der Delta-Regel



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

 $|x_1 x_2| t$ | 0 | Error | zu_addieren/subtr.

 $\Delta W(u_1,v)$ $\Delta W(u_2,v)$ $W(u_1,v)$ $W(u_2,v)$ $\Delta \theta$ 0 0 0.1 0 0 1. Epoche 100 0 1:=0+1 1:=0+1 -1:= 0-1

neue_Gewichte

0

26

Beispiel für Anwendung der Delta-Regel



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

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

Zeit $|x_1 x_2| t$ | o |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 Encelo	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	1:=1+0	0
2 Enoche	0 1	b	1	-1	0	_1	1		0	1

Delta-Regel - Beispiel

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. Epoche	0 1	0	0	0	0	0	0	0	0	0
1. Epoche	10	0	0	0	0	0	0	0	0	0
	1 1	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
	11	1,	0	1	1	1	-1	2	1,	0
1	0 0	0	0	0	0	0	0	2	1	0
3. Epoche	0 1	0	1	-1	0	-1	1	2	0	1
J. Droche	10	0	1	-1	-1	0	1	1	0	2
ernýsta naci nadk	1.1	1	0	1	. a 1 ,,	1	-1	2	1	1
100	0 0	0	0	0	0	0	0	2	1	1
4. Epoche	0 1	0	0	0	0	0	0	2	1	1
i. Epoche	10	0	1	-1	-1	0	1	1	1	2
	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
	0 0	0	0	0	0	0	0	2	1	2
6. Epoche	0 1	0	0	0	0	. 0	0	2	1	2
o. Lpoche	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

Epoche ohne Veränderung → Ende

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)

wurde hier angewandt

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.

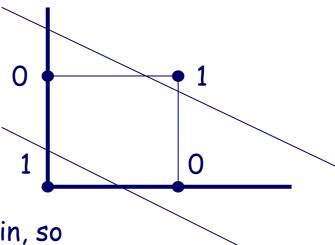
Problem: Falls das Perzeptron nicht lernt, kann nicht unterschieden werden, ob nur noch nicht genügend Schritte vollzogen wurden oder ob das Problem nicht lernbar ist. (Es gibt keine obere Schranke für die Lerndauer.)

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

Agenda



- 1. Einführung
- 2. Einfaches Perzeptron
- 3. Multi-Layer-Perzeptron
 - Vektorschreibweise der Deltaregel
 - Schichten des MLP
 - Backpropagation
 - -Overfitting
- 4. Vor- und Nachteile Neuronaler Netze

Erinnerung Perzeptron



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

Schwellwert:

 $o = \theta(x_1w_1+x_2w_2+x_3w_3+...+x_nw_n-s)$ Netz-Output:

(target) (error)

erwarteter Output: t
gemachter Fehler: e = t - o

· Lernrate:

Vektorschreibweise der Delta-Regel



Aktivität der Input-Schicht:

$$\mathbf{x} = (x_1, x_2, x_3, ..., x_n, 1)$$

Gewichtsvektor (einschl. Schwellwert):

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

Aktivität des Ausgabeneurons:

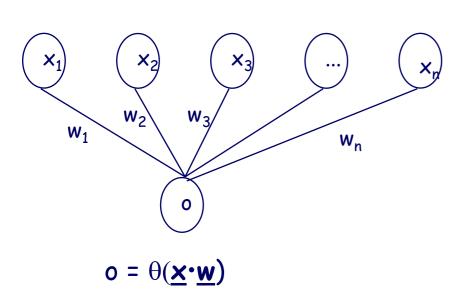
$$o = \theta(\underline{\mathbf{x}} \cdot \underline{\mathbf{w}})$$

Fehler des Ausgabeneurons († = target = erwarteter Wert): e = t-o

Gewichtsänderung:

Delta-Regel als Ableitungsregel für Perzeptron





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

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

(unter der Annahme: θ ist diff.-bar)

Fehlergradient:

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

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

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

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

$$\eta$$

Das Minus wird ab hier ignoriert, da wir den Gradienten *bergab* laufen wollen.

2-Layer-Perzeptron



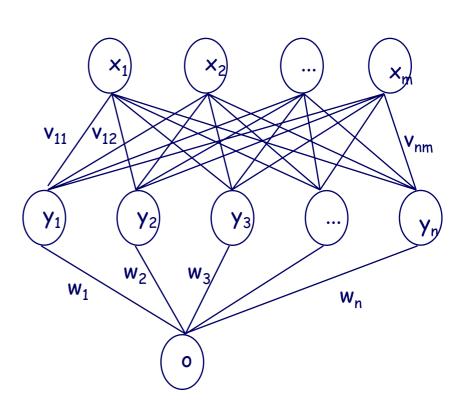
Input-Vektor
Gewichtsmatrix

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

Output o

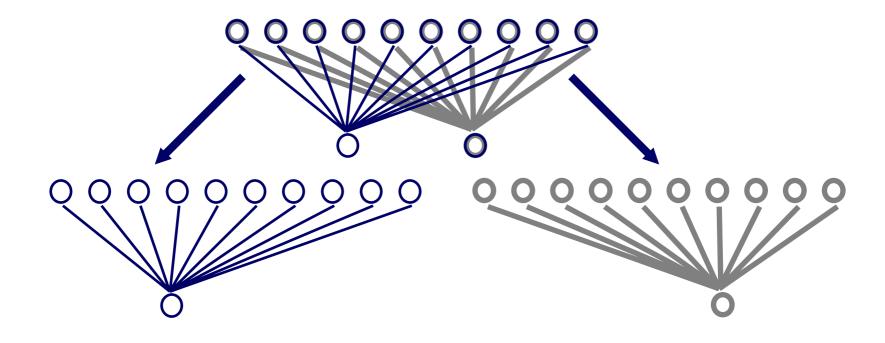
$$o = \theta(\overline{\mathbf{n}} \cdot \overline{\mathbf{x}})$$

$$\overline{\mathbf{x}} = \theta(\overline{\mathbf{n}} \cdot \overline{\mathbf{x}})$$





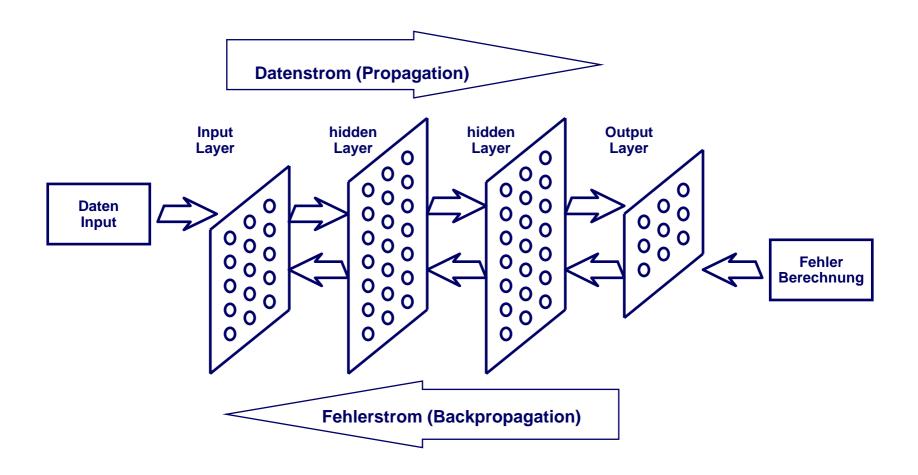
In einem 2-Schichten-Netz beeinflussen die Neuronen der 2. Schicht einander nicht, deshalb können sie voneinander getrennt betrachtet werden.



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

Backpropagation





Fehlerfunktion F (mittlerer quadratischer Fehler) für das Lernen:

$$\mathsf{F}_{\mathsf{D}} = \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$$



D Menge der Trainingsbeispiele t_d korrekter Output für d e D o_d^{α} berechneter Output für $d \in D$

Die Gewichte müssen so angepasst werden, daß der Fehler minimiert wird. Dazu bietet sich wieder das Gradientenabstiegsverfahren an. (D.h.: Bergsteigerverfahren mit Vektorraum der Gewichtsvektoren als Suchraum!)

Sei nun ein $d \in D$ gegeben. Anders geschrieben ist

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

Der Fehlergradient für wi lautet für dieses die

$$v_{11}$$
 v_{12} v_{11} v_{12} v_{13} v_{14} v_{15} v_{15} v_{17} v_{18} v_{19} v

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

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

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

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

Fehlergradient für v_{ii} 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)

$$= \partial F/\partial y_i \cdot \theta'(\underline{v}_i \cdot \underline{x}) \cdot x_i$$

$$= \partial F/\partial o \cdot \partial o/\partial y_i \cdot \theta'(\underline{v}_i \cdot \underline{x}) \cdot x_j$$

$$= \partial F/\partial o \cdot \partial \theta(\underline{\mathbf{w}} \cdot \underline{\mathbf{y}})/\partial y_i \cdot \theta'(\underline{\mathbf{v}}_i \cdot \underline{\mathbf{x}}) \cdot x_j$$

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

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

Fehler bei Info von

der Ausgabe Zwischenschicht

Info von Inputschicht



- Die schichtweise Berechnung der Gradienten ist auch für mehr als zwei Schichten möglich.
- Der Fehler wird dann 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.
- 2. Präsentiere das Muster dem Netz und berechne Output (Propagation).
- 3. Der Fehler F wird als Differenz von errechnetem und gewünschten Output ermittelt.
- 4. Die Fehlerinformationen werden durch das Netz zurückpropagiert (Backpropagation).
- 5. Die Gewichte zwischen den einzelnen Schichten werden so angepasst (Δw_{ik}), dass der mittlere Ausgabefehler für das Muster sinkt.
- 6. Abbruch, wenn Fehler auf Validierungsmenge unter gegebenem Schwellwert liegt. (Siehe spätere Folien.) Sonst gehe zu 1.

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.

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

· Es gilt bei jedem Schritt:

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

Beispiel



- · Wir wollen XOR lernen.
- Als Schwellwertfunktion verwenden wir die Sigmoid-Funktion mit Steigung 1:

$$\theta(x) = 1/(1 + e^{-x})$$

Ihre Ableitung lautet

$$\theta'(x) = e^{-x}/(1 + e^{-x})^2 = \theta(x)(1 - \theta(x))$$

- Angenommen, die Gewichte seien momentan v_{11} =0.5, v_{12} =0.75, v_{21} =0.5, v_{22} =0.5, w_{1} =0.5, w_{2} =0.5.
- Wir berechnen einmal Propagation + Backpropagation für $x_1=1, x_2=1$.

Beispiel



- Angenommen, die Gewichte seien momentan v_{11} =0.5, v_{12} =0.75, v_{21} =0.5, v_{22} =0.5, w_{1} =0.5, w_{2} =0.5.
- Wir berechnen einmal Propagation+Backpropagation für $x_1=1$, $x_2=1$.
- Propagation:
 - $y_1 := \theta(0.5 \cdot 1 + 0.75 \cdot 1) = \theta(1.25) = 0.78$
 - $y_2 := \theta(0.5 \cdot 1 + 0.5 \cdot 1) = \theta(1.0) = 0.73$
 - o:= $\theta(w_1 \cdot y_1 + w_2 \cdot y_2) = \theta(0.5 \cdot 0.78 + 0.5 \cdot 0.73) = \theta(0.755) = 0.68$
- Backpropagation Teil 1:
 - $\Delta w_i = \eta$ (t-o) y_i mit $\eta = 2 \theta'(\underline{y \cdot w})$ ergibt $\Delta w_i = 2 \theta'(\underline{y \cdot w})$ (t-o) y_i = $2 \theta'(0,755)$ (0-0.68) y_i = $2 \theta(0,755)(1-\theta(0,755))(-0.68)$ y_i = $2 \cdot 0.68(1-0.68)(-0.68)$ y_i = -0.30 y_i

$$\Delta w_1 = -0.3 \text{ y}_1 = -0.23$$
, also $w_1^{\text{neu}} := w_1^{\text{alt}} + \Delta w_1 = 0.5 - 0.23 = 0.27$
 $\Delta w_2 = -0.3 \text{ y}_2 = -0.22$, also $w_2^{\text{neu}} := w_2^{\text{alt}} + \Delta w_2 = 0.5 - 0.22 = 0.28$

Beispiel



- Backpropagation Teil 2:
 - $\Delta v_{ij} = \eta$ (t-o) x_j mit $\eta = 2 \theta'(\underline{w} \cdot \underline{y}) \cdot w_i \cdot \theta'(\underline{v}_i \cdot \underline{x})$ ergibt

```
\Delta v_{11} = 0.3 \cdot 0.5 \cdot \theta'(0.5 \cdot 1 + 0.75 \cdot 1) \cdot 0.68 \cdot 1 = 0.3 \cdot 0.5 \cdot 0.78(1 - 0.78) \cdot 0.68 \cdot 1 = -0.017
\Delta v_{12} = 0.3 \cdot 0.5 \cdot \theta'(0.5 \cdot 1 + 0.75 \cdot 1) \cdot 0.68 \cdot 1 = -0.017
\Delta v_{21} = 0.3 \cdot 0.5 \cdot \theta'(0.5 \cdot 1 + 0.5 \cdot 1) \cdot 0.68 \cdot 1 = 0.3 \cdot 0.5 \cdot 0.73(1 - 0.73) \cdot 0.68 \cdot 1 = -0.02
\Delta v_{22} = 0.3 \cdot 0.5 \cdot \theta'(0.5 \cdot 1 + 0.5 \cdot 1) \cdot 0.68 \cdot 1 = -0.02
Also
```

```
v_{11}^{\text{neu}} := v_{11}^{\text{alt}} + \Delta v_{11} = 0.5 - 0.017 = 0.483
v_{12}^{\text{neu}} := v_{12}^{\text{alt}} + \Delta v_{12} = 0.75 - 0.017 = 0.733
v_{21}^{\text{neu}} := v_{21}^{\text{alt}} + \Delta v_{21} = 0.5 - 0.02 = 0.48
v_{22}^{\text{neu}} := v_{22}^{\text{alt}} + \Delta v_{22} = 0.5 - 0.02 = 0.48
```

- Eine erneute Propagation würde nun ergeben:
 - $y_1 := \theta(0.483 \cdot 1 + 0.733 \cdot 1) = \theta(1.216) = 0.77$
 - $y_2 := \theta(0.48 \cdot 1 + 0.48 \cdot 1) = \theta(0.96) = 0.72$
 - o:= $\theta(w_1 \cdot y_1 + w_2 \cdot y_2) = \theta(0.27 \cdot 0.77 + 0.28 \cdot 0.72) = \theta(0.41) = 0.60$ statt vorher 0.68



Bemerkungen

- Jede boolesche Funktion kann durch derartiges Netz repräsentiert werden.
- Beliebige Funktionen können durch ein ANN mit drei Schichten beliebig genau approximiert werden.
- Hypothesenraum ist kontinuierlich im Gegensatz z.B. zum diskreten Hypothesenraum von Entscheidungsbaumverfahren.
- ANN kann für interne Schicht sinnvolle Repräsentationen lernen, die im vorhinein nicht bekannt sind; dies ist Gegensatz zu z.B. Inductive Logic Programming (ein Verfahren mit vorgegebenem Hintergrundwissen).

Overfitting



- Mit Overfitting beschreibt man das Problem der Überanpassung eines Modells an einen Trainings-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:
 - Trainings-Datensatz (Output bekannt)
 - ·Validierungs-Datensatz (Output bekannt)
 - Anwendungsdaten (Output unbekannt)
- (Zur Evaluierung des Verfahrens insgesamt verwendet man auch in der Anwendung Daten mit bekanntem Output (den "Testdatensatz"), um gewünschten und errechneten Output vergleichen zu können.)



Abbruchbedingungen und Überspezialisierung

- Beschreibung des Backpropagation-Algorithmus lässt Abbruchbedingung offen.
- Eine schlechte Strategie ist es, den Algorithmus solange zu iterieren, bis der Fehler für die Trainingsbeispiele unter einem vorgegebenen Schwellwert liegt, da der Algorithmus zur <u>Überspezialisierung</u> (overfitting) neigt. (Vgl. Entscheidungsbäume in KDD-Vorlesung)
- Besser: 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>: Der Fehler auf der Validierungsmenge muss nicht monoton fallen (im Gegensatz zu Fehler auf der Trainingsmenge, siehe nächste Folie)!



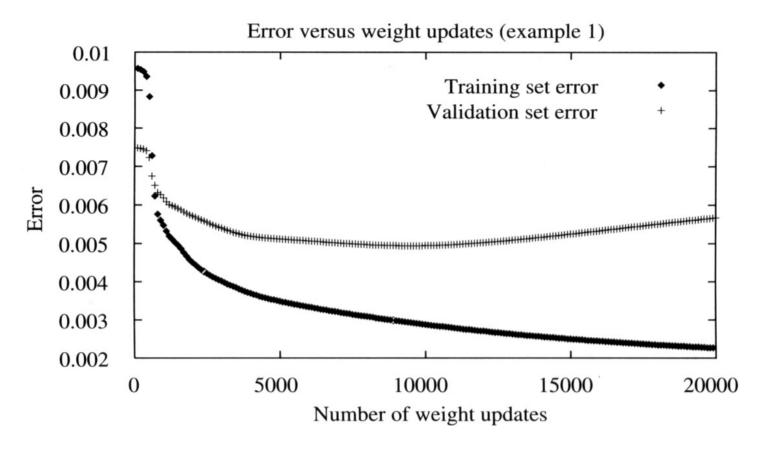


Figure 8a: Plots of error E as a function of the number of weight updates, for two different robot perception tasks (Mitchell 1997)



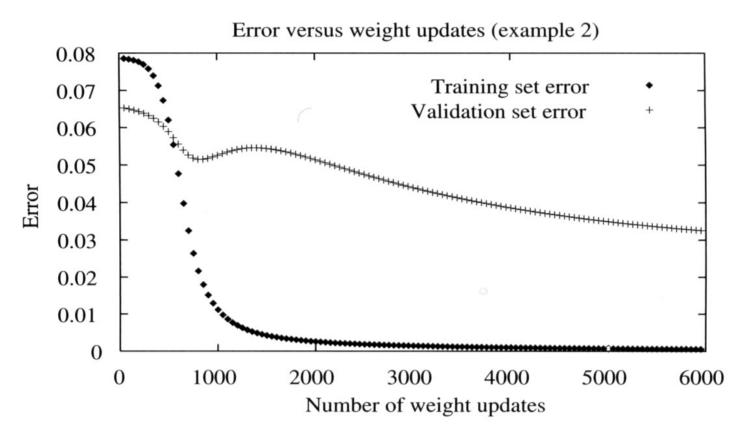


Figure 8b: Plots of error E as a function of the number of weight updates, for two different robot perception tasks (Mitchell 1997)



- Alternativ (insbes. bei wenigen Trainingsdaten) kann *k-fache Kreuzvalidierung* (*k-fold cross-validation*) 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 mittleren Fehler für die Daten aus der Validierungsmenge).
 - Der Mittelwert von i über alle k Trainingsphasen wird letztendlich verwendet, um das gegebene ANN mit allen Trainingsbeispielen zu trainieren.

Agenda

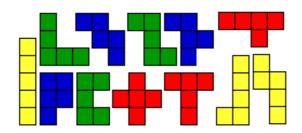


- 1. Einführung
- 2. Einfaches Perzeptron
- 3. Multi-Layer-Perzeptron
- 4. Vor- und Nachteile Neuronaler Netze

Vorteile Neuronaler Netze



· 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

Nachteile Neuronaler Netze



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