



Vorlesung Künstliche Intelligenz Wintersemester 2007/08

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.

Übersicht

1. Einführung & Grundbegriffe
 - Motivation & Definition
 - Vorbild Biologie
 - Historie der NN
 - Überblick über verschiedene Netzwerktypen
2. Einfaches Perzeptron
3. Multi-Layer-Perzeptron



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
 - Prozeßsteuerung
 - Medizin
 - Elektronische Nase
 - Schrifterkennung
 - Risikomanagement
 - Zeitreihenanalyse und -prognose
 - Robotersteuerung

Arbeitsweise Neuronale Netze

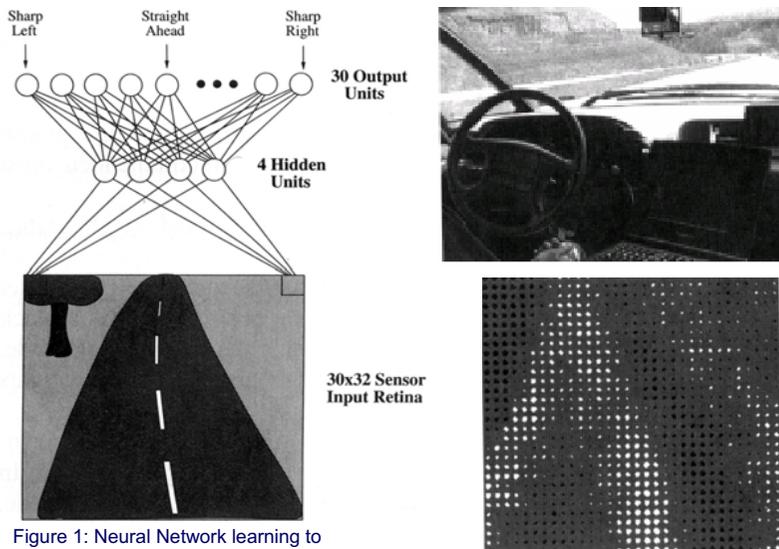


Figure 1: Neural Network learning to steer an autonomous vehicle (Mitchell 1997)

5

Charakteristische Eigenschaften von Problemen, die mit BACKPROPAGATION-ANNs gelöst werden können:

- Trainingsbeispiele sind repräsentiert durch Attribut-Wert-Paare
 - Werte können reellwertig sein
- + zu generierende Funktion (target function) kann ...
 - diskrete Funktionswerte
 - reellwertige Funktionswerte oder
 - Vektor solcher Funktionswerte
- haben
- + Trainingsbeispiele dürfen fehlerhaft sein
- lange Trainingszeiten sind akzeptabel
- + schnelle Berechnung der Funktionswerte der gelernten Funktion kann erforderlich sein
- für Menschen verständliche Interpretation der gelernten Funktion ist nicht wichtig ("Black-Box-Ansatz")

6

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.)

7

Vorbild Biologie

Gehirn



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

Geflecht aus Neuronen

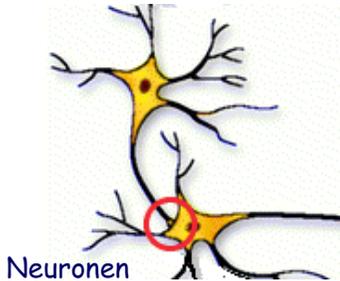


Das Gehirn besteht aus

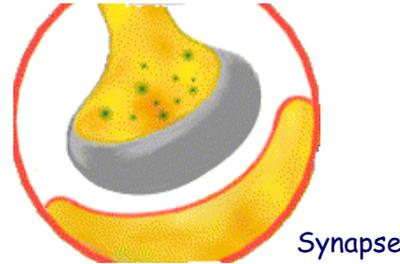
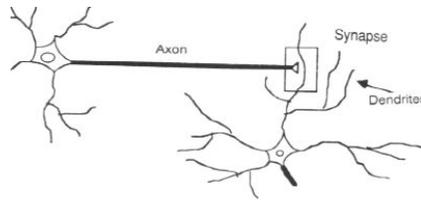
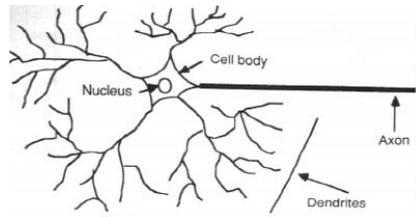
- ca. 10^{11} Neuronen, die mit
- ca. 10^4 anderen Neuronen
- durch ca. 10^{13} Synapsen verschaltet sind.

8

Vorbild Biologie



Neuronen

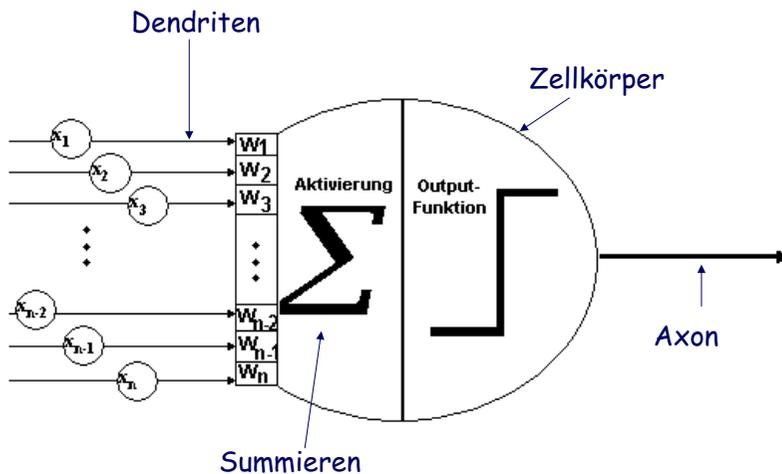


Synapse

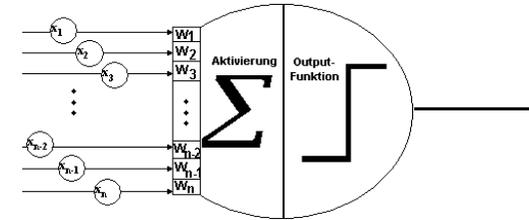
Vorbild Biologie



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, \dots, x_n werden stets mit ihren (Synapsen-)Gewichten w_1, w_2, \dots, w_n multipliziert: $a_1 = x_1 * w_1, a_2 = x_2 * w_2, \dots, a_n = x_n * w_n$.

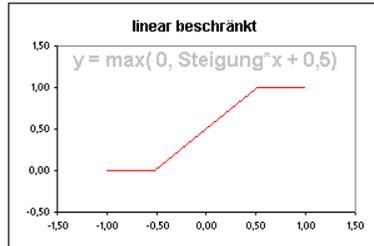
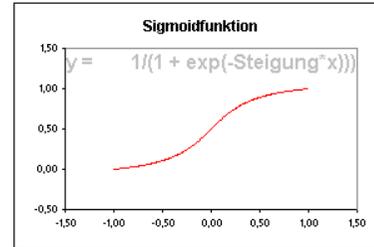
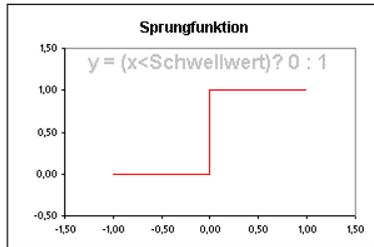
Aktivierungsfunktionen



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

1. Die am häufigsten angewandte Regel ist die **Skalarprodukt-Regel**: Die gewichteten Inputs a_1, a_2, \dots, a_n werden zur Aktivität des Neurons aufaddiert:
 $a = a_1 + a_2 + \dots + a_n$
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.

13

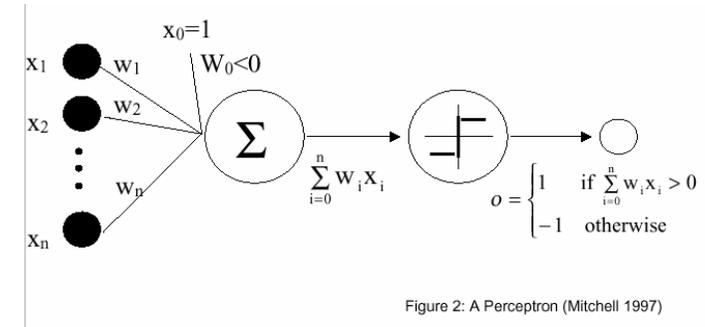
Agenda



1. Einführung
2. Einfaches Perzeptron
 - Motivation
 - Definition Perzeptron
 - Geometrische Interpretation
 - Lernen: Delta-Regel
 - XOR-Problem
3. Multi-Layer-Perzeptron

14

Perzeptron



- Input $\mathbf{x} = (x_1, \dots, x_n)$
- Gewichte $\mathbf{w} = (w_1, \dots, w_n)$
- Output (o)

15

Perzeptron

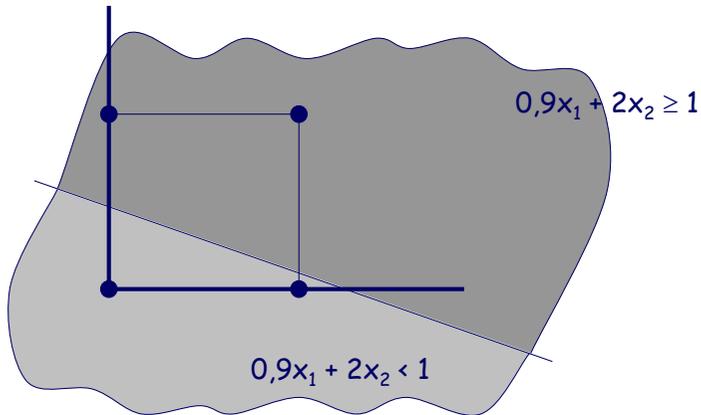


- Outputneuron hat:
 - Schwellenwert s
 - Aktivität $a := \mathbf{xw} := x_1 w_1 + \dots + x_n w_n$ (Skalarproduktregel)
 - verwendet Sprungfunktion
- Output berechnet sich wie folgt:
 - $o = 0$, falls $a < s$
 - $o = 1$, sonst.
- äquivalente Beschreibung: erweitert man die Vektoren \mathbf{x} und \mathbf{w} zu $\mathbf{y} = (x_1, \dots, x_{n+1})$ und $\mathbf{v} = (w_1, \dots, w_n, s)$, so ist
 - $o = 0$, falls $\mathbf{yv} = x_1 w_1 + \dots + x_n w_n - 1s < 0$
 - $o = 1$, sonst

16

Geometrische Interpretation

- Gleichung $xw = s$ beschreibt eine Hyperebene im n -dimensionalen Raum
- Beispiel: $0,9x_1 + 2x_2 = 1$



Konvergenz und Korrektheit der Delta-Regel

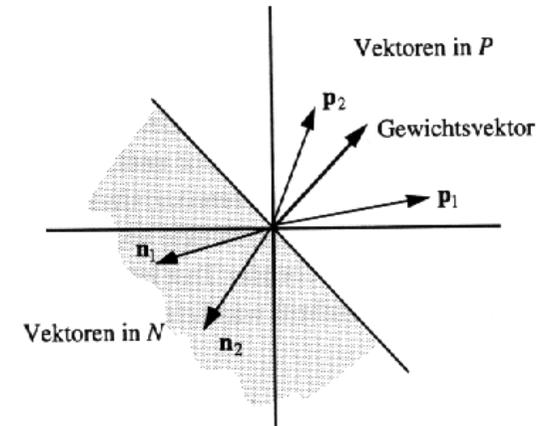


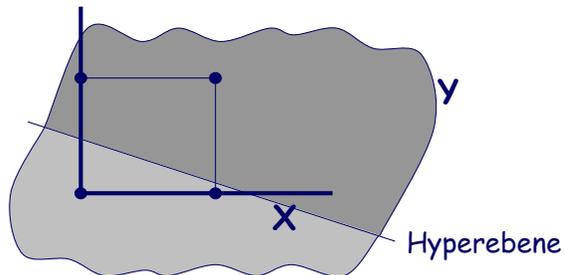
Abb. 4.8 Positiver und negativer Halbraum im Eingaberaum

entnommen Rojas, S. 84



Lernen der Gewichte

- Gegeben ist eine Menge von Beispielen.
- Überwachte Lernaufgabe: Daten sind disjunkt in zwei Mengen X, Y geteilt
- Gesucht ist der Gewichtsvektor $w (w_1, \dots, w_n)$, so dass eine Hyperebene spezifiziert wird, die die Mengen X und Y voneinander trennt.
- Mengen müssen linear trennbar sein, damit die Aufgabe lösbar ist.



Konvergenz und Korrektheit der Delta-Regel

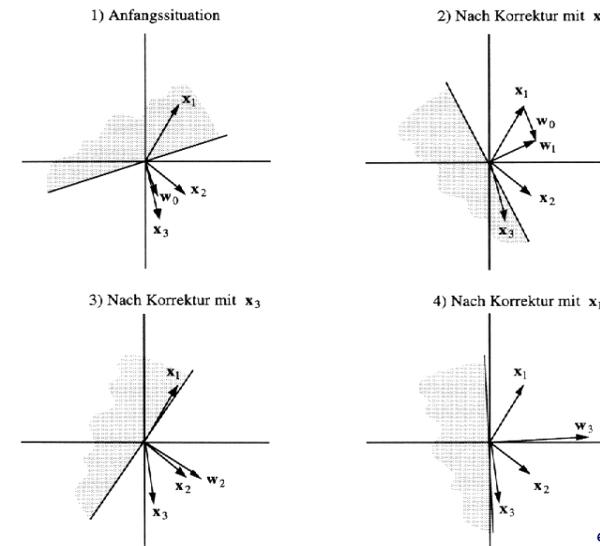


Abb. 4.9 Konvergenzverhalten des Lernalgorithmus

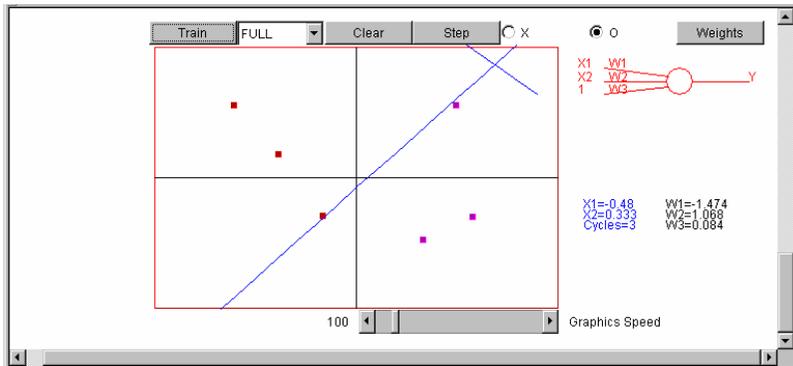
entnommen Rojas, S. 85



Delta-Regel



Menge X		Menge Y	
x ₁	x ₂	x ₁	x ₂
0,552	-0,605	0,552	0,497
0,176	-0,384	-0,304	0,579
-0,296	-0,164	-0,48	0,333



→ Demo (New38.html)

Delta-Regel



Annahme hier: - Algorithmus mit Lernrate $\eta = 1$
 - als Output nur 0 und 1 möglich (d.h. Trennung von zwei Klassen wird gelernt)

Start: Der Gewichtsvektor 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 $w_t \cdot x \leq 0$ gehe zu Addieren
 Falls $x \in Y$ und $w_t \cdot x < 0$ gehe zu Testen
 Falls $x \in Y$ und $w_t \cdot x \geq 0$ 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



- 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:

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

(mit $w_0 = -s, x_0 = 1$)

Lernrate

Beispiel für Anwendung der Delta-Regel



Wir wollen das logische Und lernen.

i	t	z
0	0	0
0	1	0
1	0	0
1	1	1

Start: Der Gewichtsvektor w_0 wird „zufällig“ generiert:
 $w_1 := 0, w_2 := 0$, Schwellenwert $\theta = w_3 := 0$

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
1. Epoche	0	0	0	0	0	0	0	0	0	0
	0	1	0	0	0	0	0	0	0	0
	1	0	0	0	0	0	0	0	0	0
	1	1	1	0	1	1	1	1	1	1
2. Epoche	0	0	0	1	-1	0	0	1	1	0
	0	1	0	1	-1	0	-1	1	0	1

In unserer Notation: w_1, w_2, w_0

Beispiel für Anwendung der Delta-Regel

Wir wollen das logische Und lernen.

i	t
0 0	0
0 1	0
1 0	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				
	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
1. Epoche	0 0	0	0	0	0	0	0	0	0	0
	0 1	0	0	0	0	0	0	0	0	0
	1 0	0	0	0	0	0	0	0	0	0
	1 1	1	0	1	1	1	-1	1	1	-1
2. Epoche	0 0	0	1	1	0	0	1	1	1	0
	0 1	0	1	-1	0	-1	1	1	0	1

$1 \times 0 + 1 \times 0 - 1_{const} \times 0 = 0$



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
1. Epoche	0 0	0	0	0	0	0	0	0	0	0
	0 1	0	0	0	0	0	0	0	0	0
	1 0	0	0	0	0	0	0	0	0	0
	1 1	1	0	1	1	1	-1	1	1	-1
2. Epoche	0 0	0	1	-1	0	0	1	1:=1+0	1:=1+0	0:= -1+1
	0 1	0	1	-1	0	-1	1	1	0	1



Beispiel für Anwendung der Delta-Regel

Addieren: 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
1. Epoche	0 0	0	0	0	0	0	0	0	0	0
	0 1	0	0	0	0	0	0	0	0	0
	1 0	0	0	0	0	0	0	0	0	0
	1 1	1	0	1	1	1	-1	1:=0+1	1:=0+1	-1:=0-1
2. Epoche	0 0	0	1	-1	0	0	1	1	1	0
	0 1	0	1	-1	0	-1	1	1	0	1



Delta-Regel - Beispiel

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
1. Epoche	0 0	0	0	0	0	0	0	0	0	0
	0 1	0	0	0	0	0	0	0	0	0
	1 0	0	0	0	0	0	0	0	0	0
	1 1	1	0	1	1	1	-1	1	1	-1
2. Epoche	0 0	0	1	-1	0	0	1	1	1	0
	0 1	0	1	-1	0	-1	1	1	0	1
	1 0	0	0	0	0	0	0	1	0	1
3. Epoche	1 1	1	0	1	1	1	-1	2	1	0
	0 0	0	0	0	0	0	0	2	1	0
	0 1	0	1	-1	0	-1	1	2	0	1
4. Epoche	1 0	0	1	-1	-1	0	1	1	0	2
	1 1	1	0	1	1	1	-1	2	1	1
	0 0	0	0	0	0	0	0	2	1	1
5. Epoche	0 1	0	0	0	0	0	0	2	1	1
	1 0	0	1	-1	-1	0	1	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
	1 0	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

29

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.

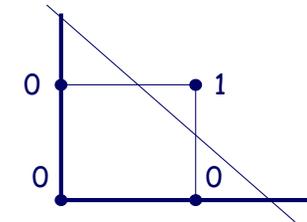
Problem: Falls das Perzeptron nicht lernt, kann nicht unterschieden werden, ob nur noch nicht genügend Schritte vollzogen wurden oder ob das Problem nicht lernbar ist. (Es gibt keine obere Schranke für die Lerndauer.)

30

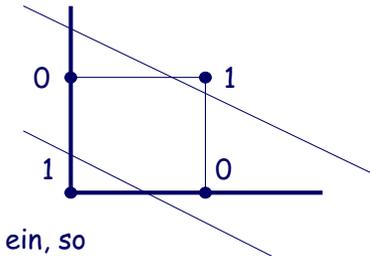
XOR-Problem



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

31

Agenda



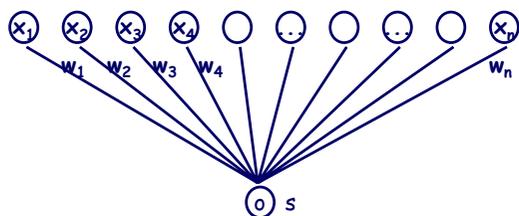
1. Einführung
2. Einfaches Perzeptron
3. Multi-Layer-Perzeptron
 - Vektorschreibweise der Deltaregel
 - Schichten des MLP
 - Backpropagation
 - Probleme der Backpropagation
 - Varianten der Backpropagation

32

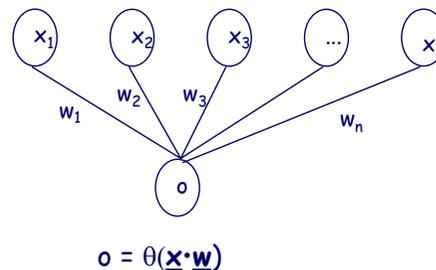
Erinnerung Perzeptron

$$\Delta w_i := \eta \cdot x_i \cdot e \quad (\text{Synapsenveränderung})$$

- Schwellwert: s
- Netz-Output: $o = \theta(x_1 w_1 + x_2 w_2 + x_3 w_3 + \dots + x_n w_n - s)$
- erwarteter Output: t (target)
- gemachter Fehler: $e = t - o$ (error)
- Lernkonstante: η



Delta-Regel als Ableitungsregel für Perzeptron



Fehlergradient:

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

$$\frac{\partial F}{\partial w_i} = \frac{\partial (o - t)}{\partial w_i}$$

$$= \frac{\partial (\theta(\underline{x} \cdot \underline{w}) - t)^2}{\partial w_i}$$

$$= 2 \underbrace{\theta'(\underline{x} \cdot \underline{w})}_{\eta} (o - t) x_i$$

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

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

(unter der Annahme: θ ist diff.-bar)



Vektorschreibweise der Delta-Regel

Aktivität der Input-Schicht:
 $\underline{x} = (x_1, x_2, x_3, \dots, x_n, 1)$

Gewichtsvektor (einschl. Schwellwert):
 $\underline{w} = (w_1, w_2, w_3, \dots, w_n, -s)$

Aktivität des Ausgabeneurons:
 $o = \theta(\underline{x} \cdot \underline{w})$

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

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



2-Layer-Perzeptron

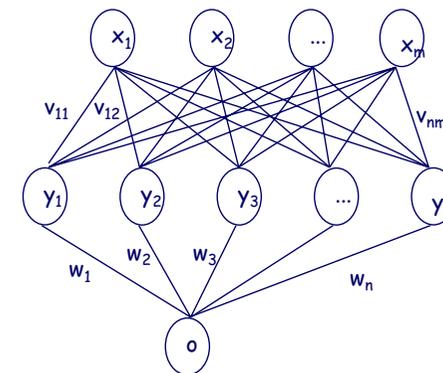
Input-Vektor \underline{x}
 Gewichtsmatrix \underline{v}

Aktivitätsvektor \underline{y}
 Gewichtsvektor \underline{w}

Output o

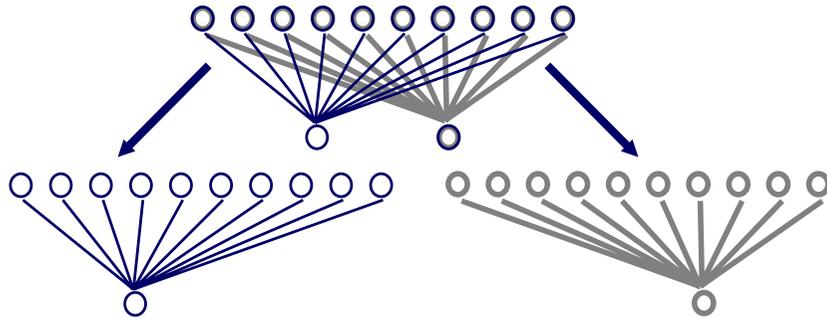
$$\underline{y} = \theta(\underline{v} \cdot \underline{x})$$

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



Multi-Layer-Perzeptron

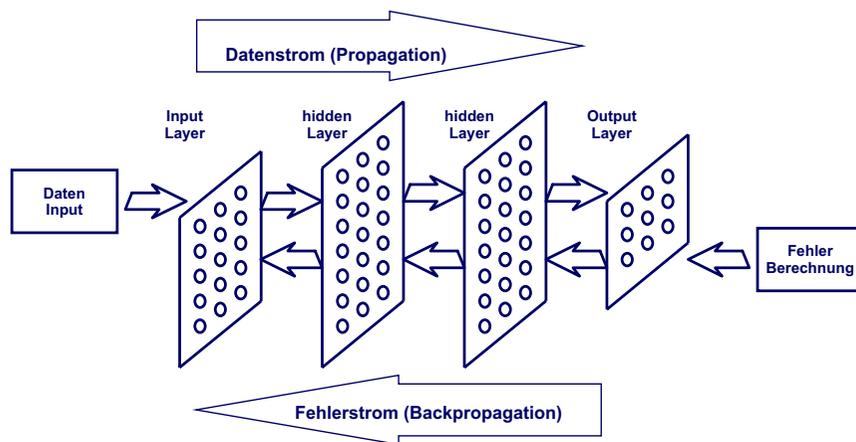
In einem 2-Schichten-Netz beeinflussen die Neuronen der 2. Schicht einander nicht, deshalb können sie voneinander getrennt betrachtet werden.



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

37

Backpropagation



38

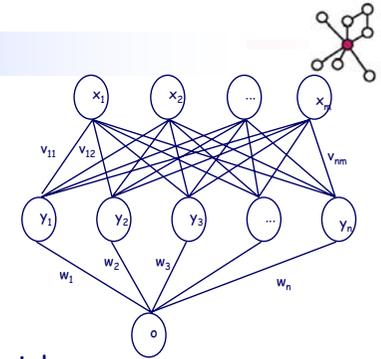
Multi-Layer-Perzeptron

Fehlerfunktion F (mittlerer quadratischer Fehler) für das Lernen:

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

wobei gilt:

D Menge der Trainingsbeispiele
 t_d korrekter Output für $d \in D$
 o_d 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!)

39

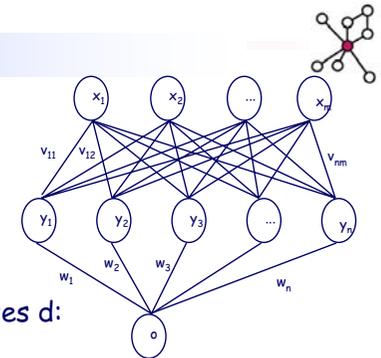
Multi-Layer-Perzeptron

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

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

Der Fehlergradient für w_i lautet für dieses d:

$$\begin{aligned} \partial F / \partial w_i &= \partial (o - t)^2 / \partial w_i \\ &= \partial (\theta(\underline{\mathbf{w}} \cdot \underline{\mathbf{y}}) - t)^2 / \partial w_i \\ &= 2 \cdot (o - t) \cdot \theta'(\underline{\mathbf{w}} \cdot \underline{\mathbf{y}}) \cdot \partial (\underline{\mathbf{w}} \cdot \underline{\mathbf{y}}) / \partial w_i \\ &= 2 \cdot (o - t) \cdot \theta'(\underline{\mathbf{w}} \cdot \underline{\mathbf{y}}) y_i \end{aligned}$$



40

Multi-Layer-Perzeptron

Fehlergradient für v_{ij} 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)

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

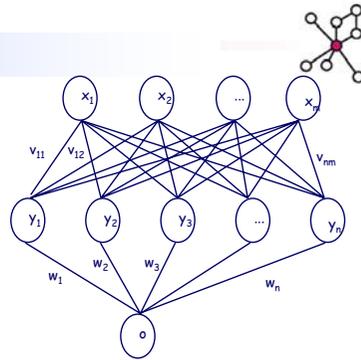
$$= \partial F / \partial o \cdot \partial o / \partial y_i \cdot \theta'(\underline{v}_i \cdot \underline{x}) \cdot x_j$$

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

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

$$= 2 \cdot (o - t) \cdot \theta'(\underline{w} \cdot \underline{y}) \cdot w_i \cdot \theta'(\underline{v}_i \cdot \underline{x}) \cdot x_j$$

Fehler bei der Ausgabe Info von Zwischenschicht Gewicht Info von Inputschicht



41

Multi-Layer-Perzeptron

- 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
 $\Delta w_{ik} = a \cdot e_k \cdot \theta'(a_k) \cdot x_i$

42

Backpropagation-Algorithmus

1. 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 gewünschtem 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.

43

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:

$$e(w + \Delta w) < e^2$$

44

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 $x_1=1$, $x_2=1$.



45

Beispiel

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

$$\begin{aligned} \Delta v_{ij} &= 2 \theta'(\mathbf{w} \cdot \mathbf{y}) w_i \cdot \theta'(\mathbf{v}_i \cdot \mathbf{x}) (o-t) x_j \\ &= 2 \theta'(0.755) w_i \cdot \theta'(0.68) 0.68 x_j \\ &= 2 \theta(0.755)(1-\theta(0.755)) w_i \cdot \theta'(0.68) 0.68 x_j \\ &= 2 \cdot 0.68(1-0.68) w_i \cdot \theta'(0.68) 0.68 x_j \\ &= 0.30 w_i \cdot \theta'(0.68) 0.68 x_j \end{aligned}$$

$$\Delta v_{11} = 0.3 \cdot 0.5 \cdot \theta'(0.5 \cdot 1 + 0.75 \cdot 1) \cdot 0.68 \cdot 1 = 0.3 \cdot 0.5 \cdot 0.78(1-0.78) \cdot 0.68 \cdot 1 = 0.017$$

$$\Delta v_{12} = 0.3 \cdot 0.5 \cdot \theta'(0.5 \cdot 1 + 0.75 \cdot 1) \cdot 0.68 \cdot 1 = 0.017$$

$$\Delta v_{21} = 0.3 \cdot 0.5 \cdot \theta'(0.5 \cdot 1 + 0.5 \cdot 1) \cdot 0.68 \cdot 1 = 0.3 \cdot 0.5 \cdot 0.73(1-0.73) \cdot 0.68 \cdot 1 = 0.02$$

$$\Delta v_{22} = 0.3 \cdot 0.5 \cdot \theta'(0.5 \cdot 1 + 0.5 \cdot 1) \cdot 0.68 \cdot 1 = 0.02$$

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



47

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:

$$\bullet y_1 := \theta(0.5 \cdot 1 + 0.75 \cdot 1) = \theta(1.25) = 0.78$$

$$\bullet y_2 := \theta(0.5 \cdot 1 + 0.5 \cdot 1) = \theta(1.0) = 0.73$$

$$\bullet o := \theta(w_1 \cdot y_1 + w_2 \cdot y_2) = \theta(0.5 \cdot 0.78 + 0.5 \cdot 0.73) = \theta(0.755) = 0.68$$

- Backpropagation Teil 1:

$$\bullet \Delta w_i = \eta (o-t) y_i \text{ mit } \eta = 2 \theta'(\mathbf{y} \cdot \mathbf{w}) \text{ ergibt}$$

$$\Delta w_i = 2 \theta'(\mathbf{y} \cdot \mathbf{w}) (o-t) y_i$$

$$= 2 \theta'(0.755) (0.68-0) y_i$$

$$= 2 \theta(0.755)(1-\theta(0.755)) 0.68 y_i$$

$$= 2 \cdot 0.68(1-0.68) 0.68 y_i$$

$$= 0.30 y_i$$

$$\Delta w_1 = 0.3 y_1 = 0.23, \text{ also } w_1^{\text{neu}} := w_1^{\text{alt}} - \Delta w_1 := 0.5 - 0.23 = 0.27$$

$$\Delta w_2 = 0.3 y_2 = 0.22, \text{ also } w_2^{\text{neu}} := w_2^{\text{alt}} - \Delta w_2 := 0.5 - 0.22 = 0.28$$



46

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. ILP (ein Verfahren mit vorgegebenem Hintergrundwissen).



48

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.)

49



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 Überspezialisierung (overfitting) neigt. (Vgl. Entscheidungsbäume in KDD-Vorlesung)
- Besser: verwende separate Menge von Beispielen zur Validierung (validation set); iteriere solange, bis Fehler für Validierungsmenge minimal ist.
- Aber: Der Fehler auf der Validierungsmenge muss nicht monoton fallen (im Gegensatz zu Fehler auf der Trainingsmenge, siehe nächste Folie)!

50

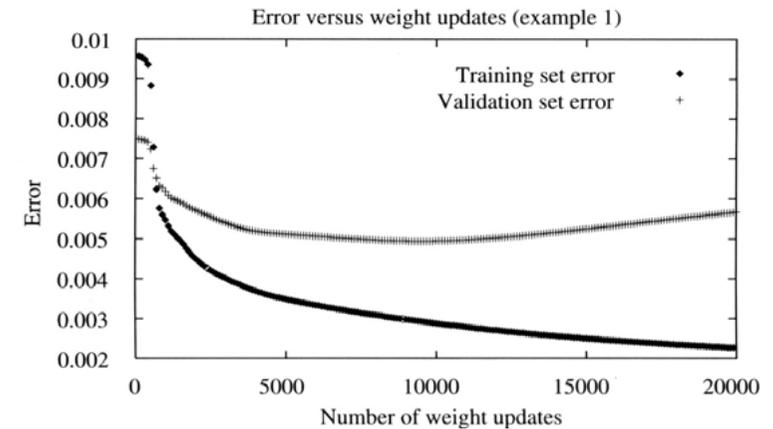


Figure 8a: Plots of error E as a function of the number of weight updates, for two different robot perception tasks (Mitchell 1997)

51

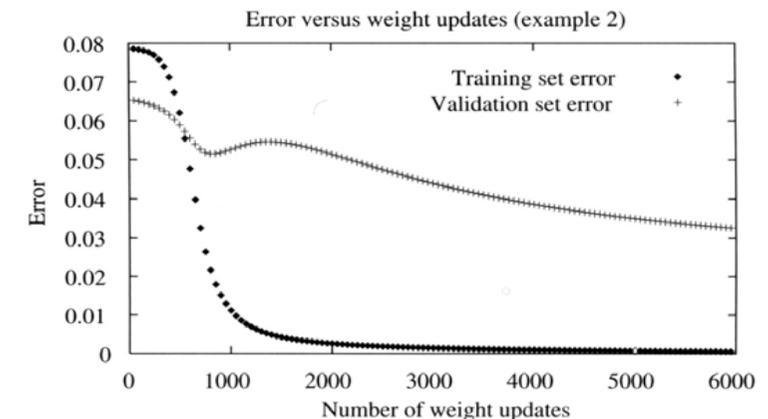


Figure 8b: Plots of error E as a function of the number of weight updates, for two different robot perception tasks (Mitchell 1997)

52



- 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.

Nachteile

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

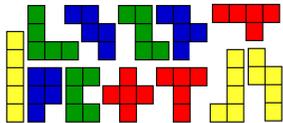
53

55



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

54