

Was bisher geschah

- ▶ biologisches Vorbild künstlicher Neuronen und künstlicher neuronaler Netze
- ▶ biologische Lernvorgänge
- ▶ mathematisches Modell: McCulloch-Pitts-Neuron
 - ▶ Boolesche Eingänge (erregend, hemmend)
 - ▶ ein Boolescher Ausgang
 - ▶ Eingangs-, Aktivierungs- und Ausgangsfunktion
 - ▶ berechnet Boolesche Funktion
 - ▶ geometrische Interpretation, Teilung des Raumes in zwei Mengen
 - ▶ linear trennbare Mengen / Boolesche Funktionen
 - ▶ Analogie zu logischen Gattern
- ▶ McCulloch-Pitts-Neuron mit (absolut) hemmenden Eingängen
- ▶ McCulloch-Pitts-Netz

Schwellwertneuronen

Idee: gewichtete Eingänge

- ▶ zur Modellierung der Stärke der synaptischen Bindung
- ▶ ermöglichen Lernen durch Änderung der Gewichte

Mathematisches Modell:

Schwellwertneuron (Perzeptron)

Eingabewerte: $x = (x_1, \dots, x_m) \in \{0, 1\}^m$

Eingangsgewichte: $w = (w_1, \dots, w_m) \in \mathbb{R}^m$

Schwellwert: $\theta \in \mathbb{R}$

Ausgabe: $a(x_1, \dots, x_m) \in \{0, 1\}$ Aktivität

Parameter eines Schwellwertneurons u :

- ▶ m_u : Anzahl der (erregenden) Eingänge
- ▶ $(w_1, \dots, w_{m_u}) \in \mathbb{R}^{m_u}$: Eingangsgewichte
- ▶ θ_u : Schwellwert

Schwellwertneuronen: Funktionen

Eingangsfunktion des Neurons u (abhängig von (w_1, \dots, w_{m_u})):

$I_u: \mathbb{R}^{m_u} \times \{0, 1\}^{m_u} \rightarrow \mathbb{R}$ mit

$$I_u(w_1, \dots, w_{m_u}, x_1, \dots, x_{m_u}) = \sum_{i=1}^{m_u} w_i x_i$$

(gewichtete Summe aller Eingänge des Neurons u)

Aktivierungsfunktion des Neurons u (abhängig von θ_u):

$A_u: \mathbb{R} \times \mathbb{R} \rightarrow \{0, 1\}$ mit

$$A_u(\theta_u, v) = \begin{cases} 1 & \text{falls } v \geq \theta_u \\ 0 & \text{sonst} \end{cases}$$

(Stufenfunktion)

Ausgabefunktion des Neurons u : $O_u: \{0, 1\} \rightarrow \{0, 1\}$ mit

$$O_u(v) = v$$

(Identität)

Schwellwertneuronen: Berechnung

vom Neuron u berechnete Funktion: $f_u: \{0, 1\}^{m_u} \rightarrow \{0, 1\}$ mit

$$\begin{aligned} f_u(x_1, \dots, x_{m_u}) &= O_u(A_u(\theta_u, I_u(w_1, \dots, w_{m_u}, x_1, \dots, x_{m_u}))) \\ &= \begin{cases} 1 & \text{falls } \langle w, x \rangle \geq \theta_u \\ 0 & \text{sonst} \end{cases} \end{aligned}$$

Wiederholung:

$\sum_{i=1}^n w_i x_i = \langle w, x \rangle$ Skalarprodukt

der Vektoren $w = (w_1, \dots, w_n)$ und $x = (x_1, \dots, x_n)$

Jedes Schwellwertneuron u mit m_u Eingängen repräsentiert eine Boolesche Funktion $f_u: \{0, 1\}^{m_u} \rightarrow \{0, 1\}$

Auch mit Schwellwertneuronen lassen sich nur linear trennbare Boolesche Funktionen berechnen (XOR nicht).

Beispiele: $\vee, \wedge, \rightarrow, ((x_1 \wedge (x_3 \vee \neg x_2)) \vee (\neg x_2 \wedge x_3))$

Schwellwertneuronen: geometrische Interpretation

Jedes Schwellwertneuron u mit m_u Eingängen teilt die Menge $\{0, 1\}^{m_u}$ der **Eingabevektoren** (Punkte im \mathbb{R}^{m_u}) in zwei Teilmengen (Teilräume des R^{m_u}):

$$\begin{aligned}f_u^{-1}(1) &= \{(x_1, \dots, x_{m_u}) \in \{0, 1\}^{m_u} \mid f(x_1, \dots, x_{m_u}) = 1\} \\ &= \{(x_1, \dots, x_{m_u}) \in \{0, 1\}^{m_u} \mid \langle w, x \rangle \geq \theta_u\}\end{aligned}$$

und

$$\begin{aligned}f_u^{-1}(0) &= \{(x_1, \dots, x_{m_u}) \in \{0, 1\}^{m_u} \mid f(x_1, \dots, x_{m_u}) = 0\} \\ &= \{(x_1, \dots, x_{m_u}) \in \{0, 1\}^{m_u} \mid \langle w, x \rangle < \theta_u\}\end{aligned}$$

Grenze: durch $\langle w, x \rangle = \theta_u$ beschriebene $(m_u - 1)$ -dimensionale Hyperebene (Teilraum)
(parallele Schnitte)

Schwellwert als Gewicht (Bias-Neuronen)

Neuron mit Schwellwert θ

Hinzufügen eines zusätzlichen Eingangs x_0 (bias neuron)
mit Wert $x_0 = 1$ (konstant)

Gewicht des Einganges x_0 : $w_0 = -\theta$

$$\begin{aligned} \sum_{i=1}^n w_i x_i \geq \theta & \quad \text{gdw.} \quad \sum_{i=1}^n w_i x_i - \theta \geq 0 \\ & \quad \text{gdw.} \quad \sum_{i=0}^n w_i x_i \geq 0 \end{aligned}$$

Überwachtes Lernen einzelner Schwellwertneuronenn

Aufgabe: Konstruktion eines Schwellwertneurons zur Berechnung einer Booleschen Funktion
 $f : \{0, 1\}^m \rightarrow \{0, 1\}$

Trainingsmenge: Menge T von Paaren (x, t) aus

- ▶ Eingabevektoren $x \in \{0, 1\}^m$ und
- ▶ Funktionswerten $t = f(x) \in \{0, 1\}$

(Werte der Funktion f an Stützstellen)

Struktur des Schwellwertneurons: Schwellwertneuron mit $m + 1$ Eingängen (bias x_0)
und Eingangsgewichten $(w_0, \dots, w_m) \in \mathbb{R}^{m+1}$

Idee: automatisches Lernen der Funktion durch (wiederholte) Änderung der Gewichte

Lernziel: Gewichte $(w'_0, \dots, w'_m) \in \mathbb{R}^{m+1}$, so dass das Schwellwertneuron die Funktion f berechnet
(Korrektheit an Stützstellen)

Δ -Regel

Idee: Lernen aus Fehlern (und deren Korrektur)

Delta-Regel:

$$\forall i \in \{0, \dots, m\} : w'_i = w_i + \Delta w_i \quad \text{mit} \quad \Delta w_i = \eta x_i (t - y)$$

- ▶ Trainingswert t
- ▶ vom Netz berechneter Wert y
- ▶ **Lernrate** $\eta \in \mathbb{R}$ (Grad der Verstärkung der Verbindung)

korrigierendes Lernen,
(falls x_i aktiv und $y \neq t$)

Beispiel: $\neg, \wedge, \rightarrow$

Δ -Lernverfahren für Schwellwertneuronen

- ▶ Beginn mit **zufälligen Eingangsgewichten** $(w_0, \dots, w_n) \in \mathbb{R}^m$ (Schwellwert als Gewicht),
- ▶ die folgenden Schritte so oft wiederholen, bis der Fehler verschwindet (oder hinreichend klein ist):
 1. Bestimmung der Schwellwertneuron-**Ausgabe** y für Trainingspaar (x, t)
 2. Bestimmung des **Fehlers** $t - y$ der tatsächlichen zur gewünschten Ausgabe vom Trainingsziel t (als Funktion $e(w_0, \dots, w_m)$ von den aktuellen Gewichten w_0, \dots, w_m),
 3. Bestimmung geeigneter **Gewichtsänderungen** Δw_i
 4. Zuordnung der **neuen Gewichte** $w'_i = w_i + \Delta w_i$ zur Verringerung des (zukünftigen) Fehlers ($e(w'_0, \dots, w'_n) < e(w_0, \dots, w_n)$)

Online-Lernen und Batch-Lernen

Lernen durch schrittweise

1. Berechnung des Fehlers
2. Berechnung der notwendigen Gewichtsänderungen
3. Änderung der Gewichte

Verfahren nach Zeitpunkt der Gewichtsänderung:

Online-Lernen Berechnung von Fehler und Gewichtsänderungen für jedes Trainingsmuster, Änderung der Gewichte sofort für jedes Trainingpaar

Batch-Lernen (Lernen in Epochen)

Epoche: Berechnung für jedes Paar der Trainingsmenge

Berechnung von Fehler und Gewichtsänderungen für die gesamte Trainingsmenge (z.B. Summe über alle Trainingpaare)

Änderung der Gewichte erst nach einer ganzen Epoche

Konvergenz des Lernverfahrens

Konvergenzsatz:

Für jede Trainingsmenge

$$\mathcal{T} \subseteq \{(x^{(i)}, t^{(i)}) \mid \forall i \in \{1, \dots, n\} : x^{(i)} \in \{0, 1\}^m \wedge t^{(i)} \in \{0, 1\}\},$$

für welche die Mengen

$$\mathcal{T}_0 = \{x \mid (x, 0) \in \mathcal{T}\} \text{ und } \mathcal{T}_1 = \{x \mid (x, 1) \in \mathcal{T}\}$$

linear trennbar sind,

terminieren sowohl Online- als auch Batch-Lernen eines Schwellwertneurons (passender Struktur) nach endlich vielen Schritten.

Die vom so trainierten Schwellwertneuron berechnete Funktion trennt die Mengen \mathcal{T}_0 und \mathcal{T}_1 voneinander.

Netze aus Schwellwertneuronen

Ein-Schicht-Schwellwertneuronen-Netz

parallele Schaltung mehrerer Schwellwertneuronen
repräsentiert Boolesche Funktionen mit mehreren
Ausgaben

Beispiel: Parallelschaltung von $x_1 \wedge x_2$ und $\neg x_1 \wedge \neg x_2$

Mehr-Schicht-Schwellwertneuronen-Netz

parallele und sequentielle Schaltung mehrerer
Schwellwertneuronen

Jede Boolesche Funktion $f : \{0, 1\}^n \rightarrow \{0, 1\}$ lässt sich durch ein
Schwellwertneuronen-Netz berechnen.

Schwellwertneuronen-Netz mit zwei Schichten genügt
(analog DNF, CNF in Aussagenlogik)

Netze aus Schwellwertneuronen

Ein-Schicht-Schwellwertneuronen-Netz

parallele Schaltung mehrerer Schwellwertneuronen
repräsentiert Boolesche Funktionen mit mehreren
Ausgaben

Beispiel: Parallelschaltung von $x_1 \wedge x_2$ und $\neg x_1 \wedge \neg x_2$

Mehr-Schicht-Schwellwertneuronen-Netz

parallele und sequentielle Schaltung mehrerer
Schwellwertneuronen

Jede Boolesche Funktion $f : \{0, 1\}^n \rightarrow \{0, 1\}$ lässt sich durch ein
Schwellwertneuronen-Netz berechnen.

Schwellwertneuronen-Netz mit zwei Schichten genügt
(analog DNF, CNF in Aussagenlogik)

Feed-Forward-Netze (FFN)

- ▶ $V = \bigcup_{k=1}^n V_k$ mit $\forall i < j \in \{1, \dots, n\} : V_i \cap V_j = \emptyset$
Zerlegung der Menge der Neuronen in n disjunkte **Schichten**
- ▶ Menge der Eingangsneuronen: V_1 (je ein Eingang)
- ▶ Menge der Ausgangsneuronen: V_n (je ein Ausgang)
- ▶ Neuronen aller anderen Schichten heißen versteckte Neuronen
- ▶ $E \subseteq \bigcup_{k=1}^{n-1} V_k \times V_{k+1}$
nur vorwärtsgerichtete Kanten zwischen benachbarten Schichten
- ▶ Gewichte bilden $m \times m$ -Matrix (mit $m = \text{Anzahl aller Neuronen}$)
- ▶ für FFN besteht die Gewichtsmatrix aus unabhängigen Blöcken
Blöcke sind die Gewichtsmatrizen zwischen den Schichten

FFN als Berechnungsmodell:

- ▶ parallele Berechnung (in den Neuronen einer Schicht)
- ▶ sequentielle Berechnung (in miteinander verbundenen Neuronen benachbarter Schichten)

Perzeptron (historisch)

1958 Frank Rosenblatt, Idee: Modell der Netzhaut (Retina)

Aufbau des Perzeptrons:

1. Schicht (Eingabeschicht) : Menge S von Stimulus-Zellen
(Verteilung)
2. Schicht (Mittelschicht) : Menge A von Assoziations-Zellen
(Vorverarbeitung)
3. Schicht (Perzeptron-Schicht) : Menge R von Response-Zellen
Muster-Assoziator aus Schwellwertneuronen
(eigentliche Verarbeitung)

Verbindungen:

- ▶ zufällig zwischen Neuronen der Eingabeschicht und Neuronen der Mittelschicht
feste Gewichte (zufällig)
- ▶ von jedem Neuron der Mittelschicht zu jedem Neuron der Ausgabeschicht
trainierbare Gewichte

Jedes Ausgabeneuron teilt die Eingabemuster in zwei Klassen
(akzeptierte und nicht-akzeptierte)

Ein-Schicht-FFN

- ▶ Abstraktion von der Eingabeschicht im historischen Perzeptron-Modell
- ▶ nur Perzeptron-Schicht (Muster-Assoziator)
- ▶ Parallele Berechnung mehrerer künstlicher Neuronen (hier Schwellwertneuronen)

Eingänge: $(x_1, \dots, x_m) \in \{0, 1\}^m$

Ausgänge: $(y_1, \dots, y_n) \in \{0, 1\}^n$

Gewichtsmatrix $W \in \mathbb{R}^{m \times n}$

Gesamtberechnung des Ein-Schicht-FFN $f : \{0, 1\}^m \rightarrow \{0, 1\}^n$ des Neurons mit gewichteter Summe als Aktivierungsfunktion:

$f(x_1, \dots, x_m) = (y_1, \dots, y_n)$ mit $\forall k \in \{1, \dots, n\} :$

$$y_k = \begin{cases} 1 & \text{falls } \sum_{i=1}^m x_i w_{ij} \geq 0 \\ 0 & \text{sonst} \end{cases}$$

(Matrixmultiplikation)

Ein-Schicht-FFN: Training mit Δ -Regel

überwachtes Lernen

Trainingsmenge: Menge von Paaren (x, t) aus

- ▶ Eingabevektoren $x \in \{0, 1\}^m$ und
- ▶ gewünschten Ausgabevektoren $t \in \{0, 1\}^n$

Lernen mit Delta-Regel für Ein-Schicht-FFN:

- ▶ Beginn mit zufälligen Eingangsgewichten $w_{ij} \in \mathbb{R}$,
- ▶ für jede Eingabe der Trainingsmenge (x, t) :
 1. Netz berechnet die Ausgabe $y = xW$,
 2. Zuordnung neuer Gewichte w'_{ij} durch Delta-Regel:

$$w'_{ij} = w_{ij} + \Delta(w_{ij}) \quad \text{mit} \quad \Delta(w_{ij}) = \eta x_i (t_j - y_j)$$

- ▶ wiederholen, bis der Fehler klein genug ist.

Das Lernverfahren mit Delta-Regel konvergiert für

- ▶ jede linear trennbare Boolesche Funktion f und
- ▶ hinreichend kleine Lernquote η

in endliche vielen Schritten zu einem Ein-Schicht-FFN, welche die Funktion f berechnet.

Künstliche Neuronen mit reellen Ein- und Ausgängen

Parameter:

Eingänge: $x_1, \dots, x_m \in \mathbb{R}^m$

Eingangsgewichte $w_1, \dots, w_m \in \mathbb{R}^m$

Ausgang: $f(\langle x, w \rangle) \in \mathbb{R}$

- ▶ Eingangsfunktion $I : \mathbb{R}^m \rightarrow \mathbb{R}$
- ▶ Aktivierungsfunktion $A : \mathbb{R} \rightarrow \mathbb{R}$
- ▶ Ausgabefunktion $O : \mathbb{R} \rightarrow \mathbb{R}$

Gesamtberechnung $f : \mathbb{R}^m \rightarrow \mathbb{R}$ des Neurons:

$$f(x_1, \dots, x_m) = O(A(I(x_1, \dots, x_m)))$$

Klassifikation durch Ein-Schicht-FFN

Klassifikation:

Zerlegung einer Menge M von Werten in (paarweise disjunkte) Klassen $\{C_1, \dots, C_n\}$, welche die Wertemenge vollständig überdecken

$$\bigcup_{i=1}^n C_i = M \quad (\forall i \neq j : C_i \cap C_j = \emptyset)$$

Klassifikation des \mathbb{R}^m durch KNN:

- ▶ Eingänge $(x_1, \dots, x_m) \in \mathbb{R}^m$
- ▶ Ausgänge $(y_1, \dots, y_n) \in \{0, 1\}^n$
für jede Klasse C_i ein Ausgabeneuron y_i
Ausgang $y_i = 1$ gdw. Eingabe $(x_1, \dots, x_m) \in C_i$

überwachtes Training des Ein-Schicht-FFN:

- ▶ zufällige Startgewichte
- ▶ schrittweise Modifikation der Gewichte zur Verringerung des Fehlers

Ein-Schicht-FFN erkennt nur linear trennbare Klassen

Problem: Wie trainiert man Mehrschicht-FFN?

Auswahl durch Mehrschicht-FFN – Beispiel

Beispiel: Auswahl aller Punkte im Einheitsquadrat

$$y = \begin{cases} 1 & \text{falls } 0 \leq x_1 \leq 1 \wedge 0 \leq x_2 \leq 1 \\ 0 & \text{sonst} \end{cases}$$

durch das 2-Schicht-FFN mit

- ▶ Eingängen x_1, x_2 und x_0 (bias)
- ▶ Ausgang y
- ▶ versteckten Neuronen z_1, \dots, z_4 und z_0 (bias)
- ▶ Gewichte der ersten Schicht (zwischen (x_0, x_1, x_2) und (z_1, \dots, z_4)):

$$W_1 = \begin{pmatrix} 1 & 0 & 1 & 0 \\ 1 & -1 & 0 & 0 \\ 0 & 0 & 1 & -1 \end{pmatrix}$$

z_1 feuert gdw. $x_1 \leq 1$, z_2 feuert gdw. $x_1 \geq 0$

z_3 feuert gdw. $x_2 \leq 1$, z_4 feuert gdw. $x_2 \geq 0$

- ▶ Gewichte der zweiten Schicht (zwischen (z_0, \dots, z_4) und y):

$$W_2 = (-7/2, 1, 1, 1, 1)^T$$

Gesamtmatrix des FFN – Beispiel

	x_0	x_1	x_2	z_0	z_1	z_2	z_3	z_4	y
x_0	0	0	0	0	1	0	1	0	0
x_1	0	0	0	0	1	-1	0	0	0
x_2	0	0	0	0	0	0	1	-1	0
z_0	0	0	0	0	0	0	0	0	-7/2
z_1	0	0	0	0	0	0	0	0	1
z_2	0	0	0	0	0	0	0	0	1
z_3	0	0	0	0	0	0	0	0	1
z_4	0	0	0	0	0	0	0	0	1
y	0	0	0	0	0	0	0	0	0

Mehr-Schicht-FFN mit linearer Aktivierung

Netzeingänge: $(x_1, \dots, x_{k_0}) \in \mathbb{R}^m$

Netzausgänge: $(y_1, \dots, y_{k_l}) \in \mathbb{R}^n$

Neuronen (l Schichten): $(z_1^0, \dots, z_{k_0}^0) \in \mathbb{R}^{k_1}$ (Eingabeneuronen)

\vdots

(versteckte Neuronen)

$(z_1^l, \dots, z_{k_l}^l) \in \mathbb{R}^{k_l}$ (Ausgabeneuronen)

Gewichtsmatrizen $W^{(j)} \in \mathbb{R}^{k_j \times k_{j+1}}$ für jedes $j \in \{0, \dots, l-1\}$

lineare Aktivierungsfunktion $I: \mathbb{R} \rightarrow \mathbb{R}$ mit $I(x) = mx$

Ausgabe des Neurons z_i^j in Schicht j :

$$f(z_1^{j-1}, \dots, z_{k_{j-1}}^{j-1}) = O(A(I(x_1, \dots, x_{k_{j-1}}))) = m \left(\sum_{l=1}^{k_{j-1}} w_{li}^{(j)} z_l^{(j-1)} \right)$$

Netzausgabe:

$$f(x_1, \dots, x_m) = m'(x_1, \dots, x_m) W^{(0)} \dots W^{(l-1)} = m'(x_1, \dots, x_m) W$$

mit $W = W^{(0)} \dots W^{(l-1)}$ (Matrixmultiplikation)

Jede Funktion, die von einem Mehr-Schicht-FFN mit linearer Aktivierung berechnet wird, kann also auch durch ein Ein-Schicht-FFN mit linearer Aktivierung berechnet werden.