

# MUSTERERKENNUNG

Vorlesung im Sommersemester 2017

Prof. E.G. Schukat-Talamazzini

Stand: 17. März 2017

## Modellneuronen und ihre Verschaltung

Feed-Forward-Netzwerke

Klassisches Perzeptron

Mehrschichtenperzeptron

MLPs für Prädiktion, Extraktion, Klassifikation

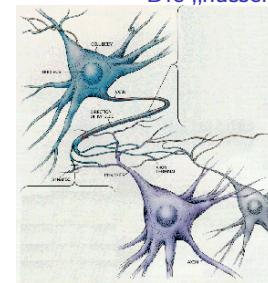
Rekurrente Netzwerke

## Teil VIII

## Künstliche Neuronale Netze

## Das menschliche Zentralnervensystem

Die „nassen“ Neuronen der anthropomorphen Blaupause



### Das Neuron

ist die elementare Verarbeitungseinheit.

- ca.  $10^{12}$  Nervenzellen
- ca.  $10^{14}$  Verbindungen

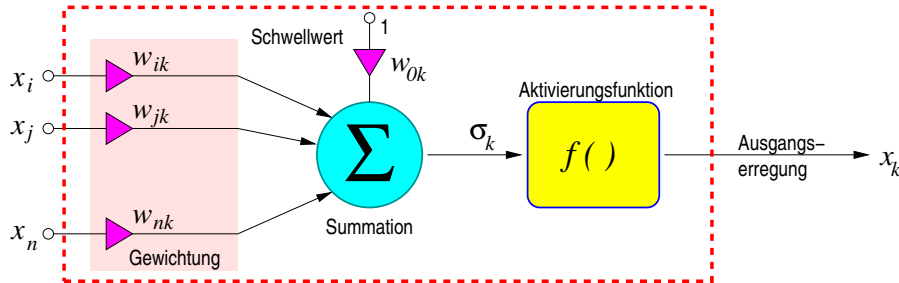
**Axon** Macht die **Grundaktivität** einer Nervenzelle in Form elektrischer Impulse verfügbar

**Synapse** Übermittlung an andere Neuronen über **synaptische** Verbindungen unterschiedlicher Stärke

**Feuern** Überschreitet die gesammelte **Eingangserregung** einen bestimmten Betrag, führt dies zu einer Aktivitätsänderung (*erhöhte Impulsfrequenz*)

## Das Modellneuron

Grundbaustein des Paradigmas *massiv-paralleler Verarbeitung*



### Parallel-Distributed Processing

Ein **Künstliches Neuronales Netz** (KNN) ist ein Arrangement von *Modellneuronen* und ihren *Verbindungen*.

### Dynamisches Verhalten

bestimmt durch

- Aktivierungsbegriff
- Kombination der Eingänge
- Aktivierungsfunktion
- Verschaltungstopologie

## Neuronales Aktivierungsniveau

Das Ausgabesignal des Modellneurons

### Ausgangserregung

Mit  $x_k(t)$  bezeichnen wir die Ausgangserregung des  $k$ -ten Modellneurons zum Zeitpunkt  $t$ .

Der **Gesamtsystemzustand** ist durch  $\mathbf{x}(t)$  charakterisiert.

#### • Wertebereich

*Feuern oder ruhen*

*Quantitatives Erregungspotenzial*

$$x_k(t) \in \{1, 0\}$$

$$x_k(t) \in \mathbb{R}$$

#### • Aktivierung

*Schwellenwert:* Neuron  $\#k$  aktiv gdw.  $x_k(t) \geq \theta$

*Feuerrate:* Neuron  $\#k$  aktiv gdw. Anzahl Spikes/Zeiteinheit  $\geq \theta$

#### • Zeitraster

dynamisch — stetige Zeitskala

dynamisch — diskrete Zeitskala

nicht-dynamisch

$$t \in \mathbb{R}$$

$$t \in \mathbb{Z}$$

$$x_k(t) \equiv x_k$$

## Neuronales Aktivierungsniveau

Die Eingangssignale des Modellneurons

### Gewichtung der Eingangserregungen

Modellierung der Stärke synaptischer Verbindungen

$$W = \begin{pmatrix} w_{11} & w_{12} & \cdots & w_{1N} \\ \vdots & & \ddots & \vdots \\ w_{N1} & w_{N2} & \cdots & w_{NN} \end{pmatrix}$$

**Exzitatorische** Synapsen

**Inhibitorische** Synapsen

$$w_{ij} > 0$$

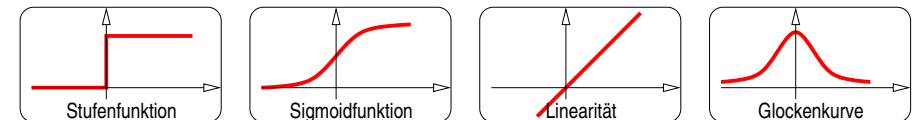
$$w_{ij} < 0$$

### Kombination der Eingangserregungen

Die gesammelte Eingangserregung des  $k$ -ten Modellneurons (zur Zeit  $t$ ) wird mit  $\sigma_k$  bzw. mit  $\sigma_k(t)$  bezeichnet und akkumuliert gemäß:

$$\sigma_k(t) = \begin{cases} \sum_j w_{jk} \cdot x_j(t) & \text{(Linearkombination)} \\ \sum_j (w_{jk} - x_j(t))^2 & \text{(Unähnlichkeit)} \end{cases}$$

## Neuronale Aktivierungsfunktionen



### Lineare Aktivierung

$$x_k = a \cdot \sigma_k + b$$

lineare Netzwerke

### Zielwertaktivierung

$$x_k = \exp \{ -C \cdot (\sigma_k - \theta_k)^2 \}$$

selbstorganisierende Karten

### Sprungfunktion

$$x_k = \begin{cases} 1 & \sigma_k > \theta_k \\ 0 & \sigma_k \leq \theta_k \end{cases}$$

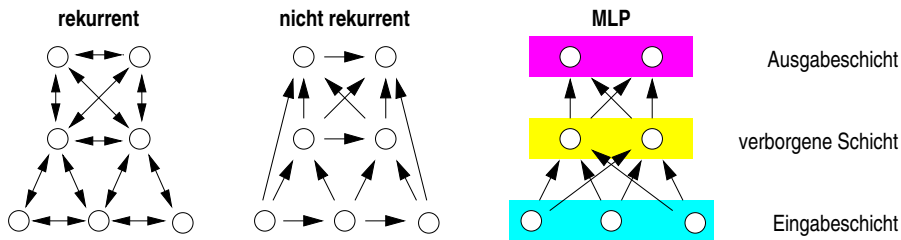
klassisches Perzeptron

### Sigmoidfunktion

$$x_k = \frac{1}{1 + \exp(-\sigma_k)}$$

Mehrschichtenperzeptron

## Neuronale Verschaltung



### Verbindungstopologie

Die Gewichtsmatrix  $\mathbf{W} \in \mathbb{R}^{N \times N}$  definiert einen **gerichteten Graphen**  $(\mathcal{N}, \mathcal{E})$  über dem Neuronenvorrat  $\mathcal{N} = \{1, \dots, N\}$  durch:

$$j \prec k \Leftrightarrow (j, k) \in \mathcal{E} \Leftrightarrow w_{jk} \neq 0, \quad (\forall j, k \in \mathcal{N})$$

### Rekurrente Netzwerke

Die Verschaltung besitzt Zyklen.

➡ *dynamisches System*

### Feed-Forward-Netzwerke

Die Verschaltung ist zyklensfrei.

➡ *endliche Berechnung*

## Feed-Forward-Netzwerke

Nichtrekurrente (zyklensfreie) künstliche neuronale Netze

### Definition

Ein KNN mit Gewichtsmatrix  $\mathbf{W} \in \mathbb{R}^{N \times N}$  und zyklensfreier neuronaler Verknüpfungsrelation ' $\prec$ ' heißt **nichtrekurrentes KNN** oder **Vorwärtsnetz**.

Die Menge  $\mathcal{N}$  aller Netzknoten zerfällt in die Teilmengen

$$\mathcal{N}_E = \{k \in \mathcal{N} \mid \neg \exists j : j \prec k\}$$

$$\mathcal{N}_A = \{j \in \mathcal{N} \mid \neg \exists k : j \prec k\}$$

$$\mathcal{N}_H = \mathcal{N} \setminus (\mathcal{N}_E \cup \mathcal{N}_A)$$

der **Eingabe**-, der **Ausgabe**- bzw. der **verborgenen** Neuronen.

### Modellneuronen und ihre Verschaltung

### Feed-Forward-Netzwerke

### Klassisches Perzeptron

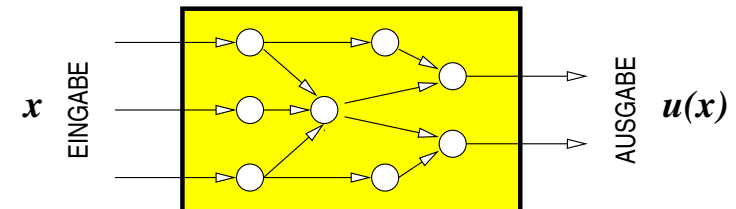
### Mehrschichtenperzeptron

### MLPs für Prädiktion, Extraktion, Klassifikation

### Rekurrente Netzwerke

## Feed-Forward-Netzwerke

Serielle Traversierbarkeit nichtrekurrenter KNNs



### Fakt

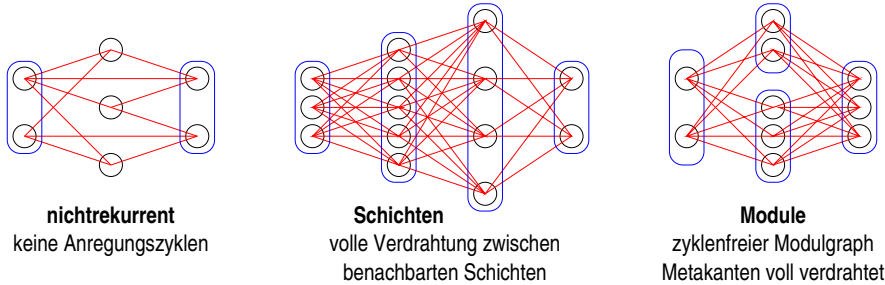
Die Knoten eines zyklensfreien KNN lassen sich in  $\prec$ -verträglicher Weise **linear anordnen**, d.h., es gilt

$$\forall j, k \in \mathcal{N} : j \prec k \Rightarrow j < k$$

nach entsprechender Umnummerierung.

## Spielarten nichtrekurrenter Netzwerke

Neuronenschichten & Verbindungsebenen



### Definition

Zerfällt die Knotenmenge  $\mathcal{N}$  eines zyklensfreien KNNs in eine Folge  $\mathcal{N}_0, \dots, \mathcal{N}_L$  disjunkter Schichten mit der Eigenschaft

$$w_{jk} \neq 0 \Leftrightarrow (\exists \ell) j \in \mathcal{N}_{\ell-1} \wedge k \in \mathcal{N}_{\ell}$$

(vollständige Verschaltung benachbarter Schichten), so sprechen wir von einem **Mehrschichtenperzeptron** (MLP = *multi-layer perceptron*).

Modellneuronen und ihre Verschaltung

Feed-Forward-Netzwerke

Klassisches Perzeptron

Mehrschichtenperzeptron

MLPs für Prädiktion, Extraktion, Klassifikation

Rekurrente Netzwerke

## Berechnungsprozeß in nichtrekurrenten Netzwerken

Zyklusfreie KNNs kommen ohne explizite *Zeitrechnung* aus

### Lemma

Die Ausgangserregungen  $x_k(t)$  der Neuronen eines Mehrschichtenperzeptrons sind nicht von der Zeit, sondern nur von den Zuständen  $\mathbf{x}_E = (x_k | k \in \mathcal{N}_E)^T$  abhängig und berechnen sich durch die Vorwärtsrekursion

$$x_k = f \left( \sum_{j \prec k} w_{jk} \cdot x_j \right) \quad \text{bzw.} \quad x_k = f \left( \sum_{j \in \mathcal{N}_{\ell-1}} w_{jk} \cdot x_j \right)$$

für alle  $\begin{cases} \text{FFN} : & k = 1, 2, \dots, N \\ \text{MLP} : & k \in \mathcal{N}_{\ell}, \ell = 1..L \end{cases}$  mit der Aktivierungsfunktion  $f(\cdot)$ .

Das Netzwerk definiert folglich eine **nichtlineare Abbildung**

$$\mathbf{u} : \begin{cases} \mathbb{R}^{|\mathcal{N}_E|} & \rightarrow \mathbb{R}^{|\mathcal{N}_A|} \\ \mathbf{x}_E & \mapsto \mathbf{u}(\mathbf{x}_E) = \mathbf{x}_A = (x_k | k \in \mathcal{N}_A)^T \end{cases}$$

## Klassisches Perzeptron

Rosenblatt 1962

### Definition

Ein einschichtiges KNN mit  $N$  Eingabeknoten,  $K$  Ausgabeknoten, der Gewichtsmatrix  $\mathbf{W} \in \mathbb{R}^{N \times K}$  und dem Systemverhalten

$$y_k(\mathbf{x}) = \text{sign} \left( \sum_{j=1}^N w_{jk} \cdot x_j \right), \quad k = 1, \dots, K$$

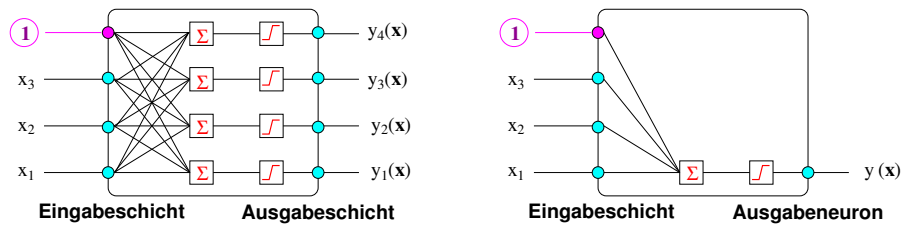
heißt **einstufiges** (oder klassisches) **Perzeptron**.

### Bemerkungen

1. Das Perzeptron kann zu Klassifikationszwecken für  $N$ -dimensionale Merkmalvektoren aus  $K$  Musterklassen herangezogen werden; Prüfgrößen sind die Aktivierungsniveaus  $y_k(\mathbf{x})$  der Ausgabeneuronen.
2. Das Verhalten der Ausgabeneuronen eines Perzeptrons ist vollständig entkoppelt; es genügt also, Perzeptren mit **einem Ausgabeneuron** zu betrachten.

# Klassisches Perzeptron mit einem Ausgabeneuron

ADALINE — *adaptive linear engine*



## Schwellenwertneuron

Feuerschwelle als  
Gewicht des 1-Neurons:

$$y = \text{sign}(\mathbf{w}^\top \mathbf{x} - \theta)$$

$$= (\theta, \mathbf{w}^\top) \cdot \begin{pmatrix} -1 \\ \mathbf{x} \end{pmatrix}$$

## Überwachtes Lernen

Ideale Zielgrößenvorgabe  $y^* \in \{-1, +1\}$   
lautet:

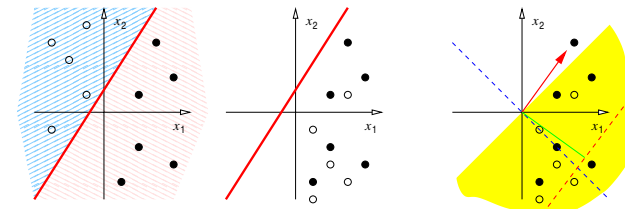
$$\text{sign}(\mathbf{w}^\top \mathbf{x}) \stackrel{!}{=} y^*$$

Mit der Hilfsgröße  $\mathbf{z} := y^* \mathbf{x}$  übersetzt in die  
Bedingung

$$\mathbf{w}^\top \mathbf{z} \stackrel{!}{\geq} 0$$

# Die Perzeptron-Lernregel

Stapelweise versus instantane Auffrischung der Gewichte



## Auffrischungsvorschrift

Korrektur von  $\mathbf{w}$  proportional zur *negativen*  
*Projektion* von  $\mathbf{w}$  auf  $\mathbf{z}$ :

$$\mathbf{w} \leftarrow \begin{cases} \mathbf{w} & \mathbf{w}^\top \mathbf{z} > 0 \\ \mathbf{w} - \eta \cdot (\mathbf{w}^\top \mathbf{z}) \cdot \mathbf{z} & \text{sonst} \end{cases}$$

Der Koeffizient  $\eta$  heißt **Lernrate**.

## Gradienten- abstieg

$$w_{jk} \leftarrow w_{jk} - \eta \Delta w_{jk}$$

**Instantan:**

$$\Delta w_{jk} = \frac{\partial \varepsilon(\mathbf{z})}{\partial w_{jk}}$$

**Stapelweise:**

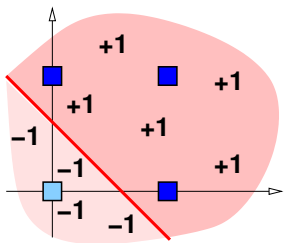
$$\Delta w_{jk} = \sum_{\mathbf{z} \in \omega} \frac{\partial \varepsilon(\mathbf{z})}{\partial w_{jk}}$$

# Was ein Perzeptron lernen kann

... und was ein Perzeptron leider **nicht** lernen kann

## Inklusive ODER-Funktion

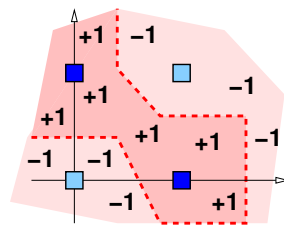
$$y(x_1, x_2) = \begin{cases} -1 & x_1 = 0 = x_2 \\ +1 & \text{sonst} \end{cases}$$



affin separierbare Muster

## Exklusive ODER-Funktion

$$y(x_1, x_2) = \begin{cases} -1 & x_1 = x_2 \\ +1 & x_1 \neq x_2 \end{cases}$$



unmöglich durch Gerade zu trennen

## Bemerkung

Nach dieser Hiobsbotschaft (1962) ruhte die Perzeptronforschung in Frieden bis zur  
(Wieder-)Entdeckung des Backpropagation-Algorithmus (1972/85).

Modellneuronen und ihre Verschaltung

Feed-Forward-Netzwerke

Klassisches Perzeptron

Mehrschichtenperzeptron

MLPs für Prädiktion, Extraktion, Klassifikation

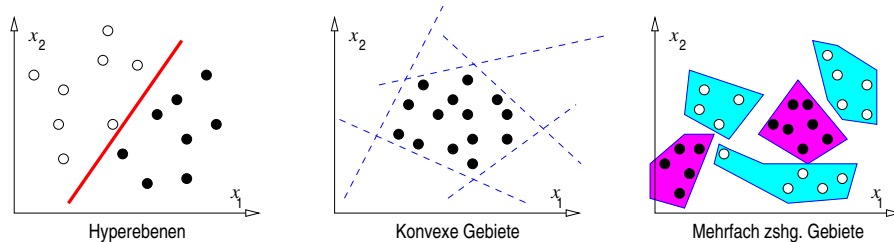
Rekurrente Netzwerke

## Das Mehrschichtenperzeptron (MLP)

Asymptotisch universeller Funktionsapproximator

### Satz

Mehrschichtennetzwerke mit nichtlinearer Aktivierungsfunktion (z.B. Signum oder Sigmoid) und mindestens zwei verborgenen Neuronenschichten ( $L \geq 3$  Stufen) können — **asymptotisch** mit wachsender Neuronenzahl — beliebige Klassengebiete bzw. Funktionen approximieren.



### Beweis.

Durch genaues Hinschauen. □

## Optimierung der neuronalen Gewichtsmatrix

### Quadratischer Vorhersagefehler

zwischen *idealer* und *neuronaler* Prüfgröße:

$$\varepsilon(\omega | \mathbf{W}) = \sum_{\mathbf{x} \in \omega} \|\mathbf{d}(\mathbf{x}) - \mathbf{y}(\mathbf{x})\|^2$$

### Stapelweise Auffrischung (*batch learning*)

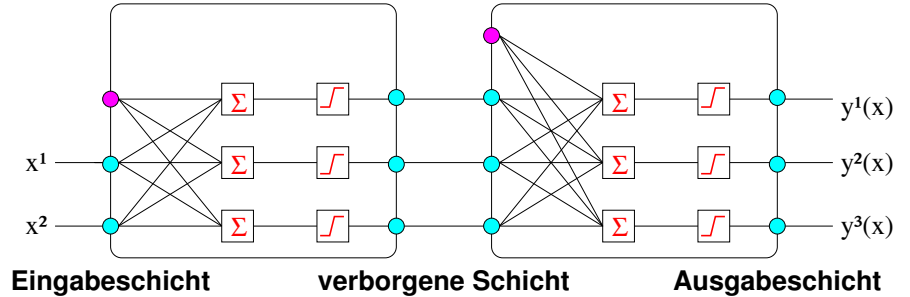
$$\mathbf{W}' = \mathbf{W} - \eta \nabla_{\mathbf{W}} \{\varepsilon(\omega | \mathbf{W})\}, \quad \nabla_{\mathbf{W}} \varepsilon = [\partial \varepsilon / \partial w_{jk}]_{j \leftarrow k}$$

### Additive Zerlegung des Gesamtfehlers

$$\frac{\partial \varepsilon(\omega | \mathbf{W})}{\partial w_{jk}} = \frac{\partial}{\partial w_{jk}} \sum_{\mathbf{x} \in \omega} \|\mathbf{d}(\mathbf{x}) - \mathbf{y}(\mathbf{x})\|^2 = \sum_{\mathbf{x} \in \omega} \frac{\partial \varepsilon(\mathbf{x} | \mathbf{W})}{\partial w_{jk}}$$

## Das Mehrschichtenperzeptron (MLP)

Beispiel: 2–3–3-Perzeptron (drei Schichten · zwei Stufen)



### Konfiguration eines MLP als Klassifikator

1.  $D$  Eingabeneuronen und  $(K - 1)$  Ausgabeneuronen  $D = 2$  und  $K = 4$  im Beispiel
2. Wieviele verborgene Schichten? eine verborgene Schicht im Beispiel
3. Wieviele Neuronen in verborgenen Schichten? 3 verborgene Neuronen im Beispiel
4. Maschinelles Lernen der Gewichtsmatrizen ( $3 \times 3$ ) und ( $4 \times 3$ ) im Beispiel

## Error-Backpropagation-Algorithmus

### Satz (Werbos 1972)

Gegeben sei ein nichtrekurrentes KNN mit Neuronen  $\mathcal{N}$ , Gewichtsmatrix  $\mathbf{W}$  und sigmoider Aktivierungsfunktion

$$f(\sigma) = 1/(1 + e^{-\sigma}) .$$

Dann gilt für die partiellen Ableitungen des quadratischen Fehlers

$$\varepsilon(\mathbf{W} | \mathbf{x}) = 1/2 \cdot \|\mathbf{d} - \mathbf{x}_A\|^2$$

zwischen der Zielvorgabe  $\mathbf{d}$  und der Netzausgabe  $\mathbf{x}_A$  zur Netzeingabe  $\mathbf{x}_E$  die Aussage

$$\partial \varepsilon / \partial w_{jk} = x_j \cdot x_k \cdot (1 - x_k) \cdot \beta_k$$

mit den rekursiv berechenbaren Größen

$$\beta_j = \frac{\partial \varepsilon}{\partial x_j} = \begin{cases} (x_j - d_j) & j \in \mathcal{N}_A \\ \sum_{k|j \prec k} \beta_k x_k (1 - x_k) w_{jk} & j \notin \mathcal{N}_A \end{cases} .$$

## Beweis.

Unsere gewählte Aktivierungsfunktion besitzt eine besonders schöne erste Ableitung:

$$f'(\sigma) = f(\sigma) \cdot (1 - f(\sigma))$$

Für die gesuchte Gewichtableitung ergibt die Kettenregel

$$\frac{\partial \varepsilon}{\partial w_{ij}} = \frac{\partial \varepsilon}{\partial x_j} \cdot \frac{\partial x_j}{\partial w_{ij}}$$

Der **rechte Faktor** wird wie folgt zerlegt:

$$\frac{\partial x_j}{\partial w_{ij}} = \frac{\partial x_j}{\partial \sigma_j} \cdot \frac{\partial \sigma_j}{\partial w_{ij}}$$

Wegen  $\frac{\partial x_j}{\partial \sigma_j} = \frac{\partial f(\sigma_j)}{\partial \sigma_j} = f'(\sigma_j) = f(\sigma_j) \cdot (1 - f(\sigma_j)) = x_j(1 - x_j)$

und

$$\frac{\partial \sigma_j}{\partial w_{ij}} = \frac{\partial}{\partial w_{ij}} \sum_{k|k \prec j} x_k w_{kj} = x_i$$

folgt unmittelbar

$$\frac{\partial x_j}{\partial w_{ij}} = f'(\sigma_j) \cdot x_i = x_j(1 - x_j) \cdot x_i$$

□

Neuronale Architektur   Feed-Forward   SLP   **MLP**   Anwendung   Rekurrente Netzwerke

## Error-Backpropagation-Algorithmus

(Algorithmus)

- 1 Initialisiere Gewichtsmatrix  $\mathbf{W}$  und setze  $\nu = 0$ .
- 2 Initialisiere Akkumulatoren  $\Delta^{(\nu)} w_{ij} = 0$ .
- 3 Für alle Klassen  $\kappa = 1, \dots, K$  und für alle Muster  $\mathbf{x}_E \in \omega_\kappa$ :
  - a Berechne Eingangserregungen  $\boldsymbol{\sigma}(\mathbf{x}_E)$ .  
Berechne Ausgangserregungen  $\mathbf{x}(\mathbf{x}_E)$ .
  - b Für alle Ausgabeneuronen  $j \in \mathcal{N}_A$ :

$$\beta_j = x_j - \begin{cases} 1 & \text{Neuron } j \text{ codiert } \Omega_\kappa \\ 0 & \text{sonst} \end{cases}$$

- c Für alle anderen Neuronen  $j \in \mathcal{N} \setminus \mathcal{N}_A$ :

$$\beta_j = \sum_{k|j \prec k} \beta_k \cdot f'(\sigma_k) \cdot w_{jk}$$

- d Für alle Neuronen  $i, j \in \mathcal{N}$  mit  $i \prec j$ :

$$\Delta^{(\nu)} w_{ij} \leftarrow \Delta^{(\nu)} w_{ij} + \beta_j \cdot f'(\sigma_j) \cdot x_i$$

- 4 Berechnen der neuen Gewichte  $\mathbf{W} \leftarrow \mathbf{W} + \eta \cdot \Delta^{(\nu)} \mathbf{W}$

- 5 Abbruch? Sonst  $\nu \leftarrow \nu + 1$  und weiter bei  $\leadsto$  2.

(summiert über  $\mathcal{A}$ )

## Beweis.

Der **linke Faktor** beschreibt die Fehlerkomponente des  $j$ -ten Neurons. Für Ausgabeneuronen gilt einfach

$$\beta_j = \frac{\partial \varepsilon}{\partial x_j} = \frac{\partial}{\partial x_j} \frac{1}{2} \cdot \sum_{k \in \mathcal{N}_A} (d_k - x_k)^2 = (x_j - d_j)$$

Für alle anderen Neuronen gilt die Rekursionsformel

$$\begin{aligned} \beta_j = \frac{\partial \varepsilon}{\partial x_j} &= \sum_{k|j \prec k} \frac{\partial \varepsilon}{\partial x_k} \cdot \frac{\partial x_k}{\partial x_j} \\ &= \sum_{k|j \prec k} \frac{\partial \varepsilon}{\partial x_k} \cdot \frac{\partial x_k}{\partial \sigma_k} \cdot \frac{\partial \sigma_k}{\partial x_j} \\ &= \sum_{k|j \prec k} \beta_k \cdot f'(\sigma_k) \cdot w_{jk} \end{aligned}$$

Insgesamt ergibt sich mit diesen Größen der Ausdruck

$$\frac{\partial \varepsilon}{\partial w_{ij}} = \beta_j \cdot f'(\sigma_j) \cdot x_i$$

für die partiellen Gewichtableitungen. Die Behauptung des Satzes folgt durch Einsetzen der Sigmoidfunktion. □

Neuronale Architektur   Feed-Forward   SLP   **MLP**   Anwendung   Rekurrente Netzwerke

## Error-Backpropagation-Algorithmus

Mächtiges Lernverfahren mit mangelhaften Konvergenzeigenschaften

### Anwendbarkeit

Alle nichtrekurrenten Netzwerke (incl. MLPs und TDNNs)

Alle differenzierbaren Aktivierungsfunktionen

Alle differenzierbaren Eingabekombinationen (incl.  $\mathbf{w}_k^\top \mathbf{x}$  und  $\|\mathbf{w}_k - \mathbf{x}\|$ )

### Konvergenzverhalten

Gradientenabstiegsverfahren

findet nur lokale Fehlerminima

Heuristische Schrittweite (Lernrate  $\eta$ )   ⬆️ Oszillation   ⬇️ Konvergenzrate

### Konvergenzbeschleunigung

Der Korrekturterm wird durch **Momentumterm** erweitert:

$$\Delta^{(\nu)} w_{ij} = \eta \cdot \frac{\partial \varepsilon}{\partial w_{ij}} + \alpha \cdot \Delta^{(\nu-1)} w_{ij}$$

Resilienter EBP-Algorithmus, Gauß-Newton, Levenberg-Marquardt, ...

## Mehrschichtenperzeptron

Serielle Notation mit abstrakten Verarbeitungsstufen



$1 + L$  Variablenschichten

zu je  $D_\ell$  Komponenten:

$$\Omega_\ell \cong \mathbb{R}^{D_\ell}$$

$L$  Verarbeitungsstufen

zu je  $D_\ell$  Einzelfunktionen:

$$f_{\ell,j} : \Omega_{\ell-1} \rightarrow \mathbb{R}$$

für  $j = 1, \dots, D_\ell$

Beispielverarbeitungsstufen

Lineargewichtungsstufe

$$f_\ell : \begin{cases} \mathbb{R}^n & \rightarrow \mathbb{R}^m \\ \mathbf{x} & \mapsto \mathbf{A}^\top \cdot \mathbf{x} + \mathbf{b} \end{cases}$$

Kompandierungsstufe (z.B. Sigmoid)

$$f_\ell : \begin{cases} \mathbb{R}^n & \rightarrow \mathbb{R}^n \\ \mathbf{x} & \mapsto (\varphi(x_1), \dots, \varphi(x_n))^\top \end{cases}$$

Softmaxstufe o.ä.

$$f_\ell : \begin{cases} \mathbb{R}^n & \rightarrow \mathbb{R}^k \\ \mathbf{x} & \mapsto (e^{x_j} / \sum_i e^{x_i})_{j=1..k} \end{cases}$$

## Kosten, Lerndaten und Gütefunktion

Optimierung der freien MLP-Parameter  $\theta$

Kostenfunktional

$\varepsilon(\hat{\mathbf{y}}, \mathbf{y}^*)$  für

Wunschergebnis  $\mathbf{y}^*$  und

MLP-Ausgabe  $\hat{\mathbf{y}} = f_{1:L}(\mathbf{x})$

Beispiel für

Regression/Klassifikation

$$\varepsilon(\hat{\mathbf{y}}, \mathbf{y}^*) = \frac{1}{2} \|\hat{\mathbf{y}} - \mathbf{y}^*\|^2 \quad (\text{OLS})$$

$$\varepsilon(\hat{\mathbf{y}}, k^*) = -\log(\hat{y}_{k^*}) \quad (\text{MAP})$$

Stichprobe  $\omega = (\mathbf{x}^{(t)}, \mathbf{y}^{(t)})_{t=1..T}$  für überwachtes Lernen

$$\text{Näherungsziel: } \mathcal{E}[\varepsilon(f_\theta(\mathbb{X}), \mathbb{Y})] \approx \sum_{t=1}^T \varepsilon(f_\theta(\mathbf{x}^{(t)}), \mathbf{y}^{(t)}) \quad (\text{LLN})$$

Gradientenabstieg

$$\frac{\partial \mathcal{E}[\varepsilon(f_\theta(\mathbb{X}), \mathbb{Y})]}{\partial \theta} \approx \sum_{t=1}^T \frac{\partial \varepsilon(f_\theta(\mathbf{x}^{(t)}), \mathbf{y}^{(t)})}{\partial \theta}$$

## Berechnung der Vorwärtswariablen

Schichtenparallel von Stufe zu Stufe



Eingabe  $\rightsquigarrow$  Ausgabe

Zweistufen-MLP zur Klassifikation

$$\begin{aligned} \mathbf{x}_0 &:= \mathbf{x} \in \mathbb{R}^{D_0} \\ \mathbf{x}_1 &:= f_1(\mathbf{x}_0) \\ \mathbf{x}_2 &:= f_2(\mathbf{x}_1) \\ \dots &\dots \\ \mathbf{x}_\ell &:= f_\ell(\mathbf{x}_{\ell-1}) \\ \dots &\dots \\ \mathbf{x}_L &:= f_L(\mathbf{x}_{L-1}) \end{aligned}$$

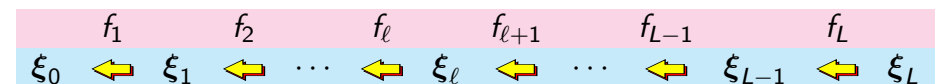
$$\begin{aligned} x_{0,i} &= x_i && \text{für } i = 1..N \\ x_{1,j} &= W_{0j} + \sum_i W_{ij} \cdot x_{0,i} && \text{für } j = 1..M \\ x_{2,j} &= (1 + e^{-x_{1,j}})^{-1} && \text{für } j = 1..M \\ x_{3,k} &= V_{0k} + \sum_j V_{jk} \cdot x_{2,j} && \text{für } k = 1..K \\ x_{4,k} &= (1 + e^{-x_{3,k}})^{-1} && \text{für } k = 1..K \\ \hat{y}_k &= x_{4,k} \end{aligned}$$

$$f_{1:L}(\mathbf{x}) = (f_1 \circ f_2 \circ \dots \circ f_L)(\mathbf{x}) = f_L(f_{L-1}(\dots(f_2(f_1(\mathbf{x})))\dots))$$

## Ableitung nach den Vorwärtswariablen

Fehlerrückführung vom MLP-Ausgang zum MLP-Eingang

$$\xi_{\ell,j} \stackrel{\text{def}}{=} \partial \varepsilon / \partial x_{\ell,j} \quad \text{für alle } \left\{ \begin{array}{l} \ell = 0..L \\ j = 1..D_\ell \end{array} \right\}$$



Rekursionsbeginn

$$\xi_{L,j} = \frac{\partial}{\partial x_{L,j}} \varepsilon(\mathbf{x}_L, \mathbf{y}^*)$$

OLS-Regression

$$\begin{aligned} \xi_{L,j} &= \frac{\partial}{\partial x_{L,j}} \frac{1}{2} \sum_i (x_{L,i} - y_i^*)^2 \\ &= x_{L,i} - y_i^* \end{aligned}$$

Rekursionsschritt

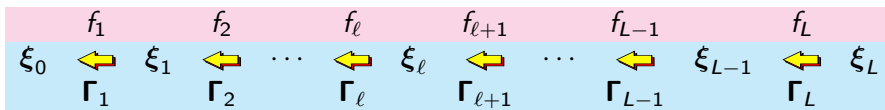
$$\begin{aligned} \xi_{\ell-1,j} &= \frac{\partial \varepsilon}{\partial x_{\ell-1,j}} \\ &= \sum_k \frac{\partial \varepsilon}{\partial x_{\ell,k}} \cdot \frac{\partial x_{\ell,k}}{\partial x_{\ell-1,j}} \\ &= \sum_k \xi_{\ell,k} \cdot \gamma_{\ell,j,k} \end{aligned}$$



## Lokale Ableitungen

zwischen den Variablen zweier benachbarter MLP-Stufen

$$\gamma_{\ell,j,k} \stackrel{\text{def}}{=} \partial x_{\ell,k} / \partial x_{\ell-1,j} \quad \text{für alle } \begin{cases} \ell = 1 \dots L \\ j = 1 \dots D_{\ell-1} \\ k = 1 \dots D_{\ell} \end{cases}$$



### Aktivierungsstufe

$$\gamma_{\ell,j,k} = \frac{\partial}{\partial x_{\ell-1,j}} \varphi(x_{\ell-1,k}) = \begin{cases} \varphi'(x_{\ell-1,k}) & j = k \\ 0 & j \neq k \end{cases}$$

### Linearkombination

$$\gamma_{\ell,j,k} = \frac{\partial}{\partial x_{\ell-1,j}} \left( b_{\ell,k} + \sum_i A_{\ell,i,k} \cdot x_{\ell-1,i} \right) = A_{\ell,j,k}$$

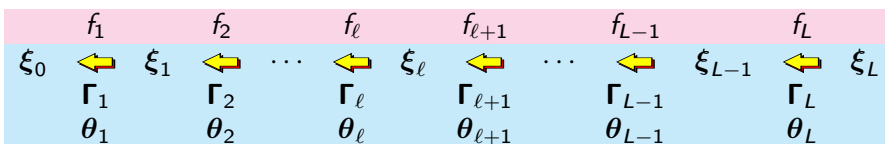
### Softmaxstufe

$$\gamma_{\ell,j,k} = \frac{\partial}{\partial x_{\ell-1,j}} \frac{e^{x_{\ell-1,k}}}{\sum_i e^{x_{\ell-1,i}}} = \begin{cases} -x_{\ell,j} \cdot x_{\ell,k} + x_{\ell,k} & j = k \\ -x_{\ell,j} \cdot x_{\ell,k} + 0 & j \neq k \end{cases}$$

## Ableitung nach den MLP-Parametern

„Synaptische Gewichte“ des Neuronennetzes · Kompandierungsparameter

$$\frac{\partial \varepsilon}{\partial \theta_{\ell,i}} = \sum_k \frac{\partial \varepsilon}{\partial x_{\ell,k}} \cdot \frac{\partial x_{\ell,k}}{\partial \theta_{\ell,i}} = \sum_k \xi_{\ell,k} \cdot \frac{\partial}{\partial \theta_{\ell,i}} f_{\ell,k}(x_{\ell-1}) \quad \text{für alle } \begin{cases} \ell = 1 \dots L \\ i = 1 \dots I_{\ell} \end{cases}$$



### Linearkoeffizienten

$$\frac{\partial \varepsilon}{\partial A_{\ell,j,k}} = \frac{\partial \varepsilon}{\partial x_{\ell,k}} \cdot \frac{\partial x_{\ell,k}}{\partial A_{\ell,j,k}} = \xi_{\ell,k} \cdot \frac{\partial}{\partial A_{\ell,j,k}} \left( b_{\ell,k} + \sum_i A_{\ell,i,k} \cdot x_{\ell-1,i} \right) = \xi_{\ell,k} \cdot x_{\ell-1,j}$$

### Aktivierungsschwellen

$$\frac{\partial \varepsilon}{\partial b_{\ell,k}} = \frac{\partial \varepsilon}{\partial x_{\ell,k}} \cdot \frac{\partial x_{\ell,k}}{\partial b_{\ell,k}} = \xi_{\ell,k} \cdot \frac{\partial}{\partial b_{\ell,k}} \left( b_{\ell,k} + \sum_i A_{\ell,i,k} \cdot x_{\ell-1,i} \right) = \xi_{\ell,k} \cdot 1$$

## Kompandierungsfunktionen und ihre Ableitungen

nichtlineare, monotone Skalentransformationen für Aktivierungsstufen

Name	Funktion	Ableitung
Identität <sup>1</sup>	id(x) x	1
Stufe <sup>2</sup>	sign(x) 2 · I <sub>x&gt;0</sub> (x) - 1	0 für x ≠ 0
Sigmoid <sup>3</sup>	σ(x) 1 / (1 + e <sup>-x</sup> )	σ(x) · (1 - σ(x))
Tangens <sup>4</sup>	tanh(x) (e <sup>2x</sup> - 1) / (e <sup>2x</sup> + 1)	1 - tanh <sup>2</sup> (x)
Ellbogen <sup>5</sup>	φ <sub>+</sub> (x) max(0, x)	I <sub>x&gt;0</sub> (x)
μ-Gesetz <sup>1</sup>	φ <sub>μ</sub> (x) sign(x) · log(1 + μ x )	( x  + 1/μ) <sup>-1</sup>
μ-invers <sup>1</sup>	φ <sub>μ</sub> (x) sign(x) · (e <sup> x </sup> - 1)/μ	e <sup> x </sup> /μ

### Bemerkung

Wertebereiche: <sup>1</sup> ℝ <sup>2</sup> {-1, 0, 1} <sup>3</sup> (0, 1) <sup>4</sup> (-1, 1) <sup>5</sup> ℝ<sub>+</sub>

## Fehlerrückführung in Matrixschreibweise

L-Stufen-MLP · Regressionsaufgabe · Methode der kleinsten Quadrate

(Algorithmus)

- 1 Initialisiere die Parameterfelder  $\mathbf{A}_{\ell}$  und  $\mathbf{b}_{\ell}$  für  $\ell = 1 \dots L$
- 2 Initialisiere Akkumulatorfelder  $\Delta_{\ell}^A$  und  $\Delta_{\ell}^b$  zu Null
- 3 Für alle Lernstichprobenpaare  $(\mathbf{x}, \mathbf{y}) \in \omega$  berechne ...

- a Vorwärtsvariable  $\mathbf{x}_{\ell} = \begin{cases} \mathbf{x} & \ell = 0 \\ \mathbf{A}_{\ell}^{\top} \cdot \mathbf{x}_{\ell-1} + \mathbf{b}_{\ell} & \text{linear} \\ \varphi(\mathbf{x}_{\ell-1}) & \text{skalar} \end{cases}$
- b Rückwärtsvariable  $\xi_{\ell} = \begin{cases} \mathbf{x}_L - \mathbf{y} & \ell = L \\ \mathbf{A}_{\ell+1} \cdot \xi_{\ell+1} & \text{linear} \\ \text{diag}(\varphi'(\mathbf{x}_{\ell})) \cdot \xi_{\ell+1} & \text{skalar} \end{cases}$
- c Aktualisiere  $\Delta_{\ell}^A$  und  $\Delta_{\ell}^b$  um  $\mathbf{x}_{\ell} \cdot \xi_{\ell}^{\top}$  bzw.  $\xi_{\ell}$

- 4 Verschiebe die alten Parameter in Gradientenrichtung

- 5 Abbruch oder weiter bei  $\rightsquigarrow$  2

(zumfliegIA)

## Modellneuronen und ihre Verschaltung

## Feed-Forward-Netzwerke

## Klassisches Perzeptron

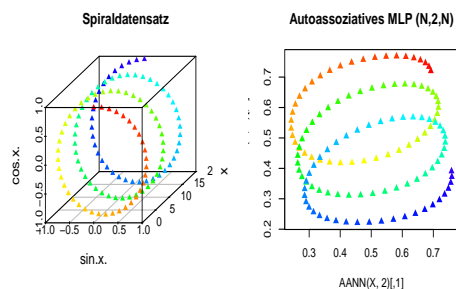
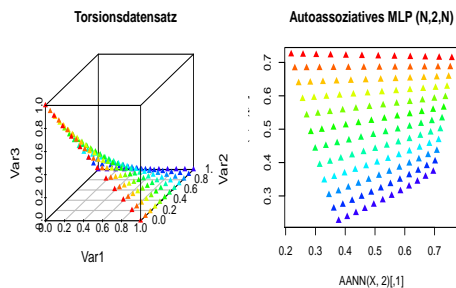
## Mehrschichtenperzeptron

## MLPs für Prädiktion, Extraktion, Klassifikation

## Rekurrente Netzwerke

## Autoassoziatives Mehrschichtenperzeptron

```
neuralnet::neuralnet (formula, data, hidden=1, algorithm='rprop+', ...)
```



## Funktionsargumente

```
formula Output/Input-Variablen
data Datensatz
hidden Zwischenschichtengrößen
algorithm {backprop, sag, slr
           rprop+, rprop-}
return Objekt der Klasse nn
```

## R-Programmcode

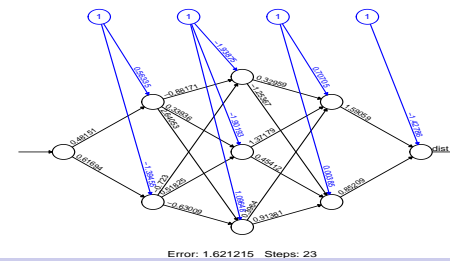
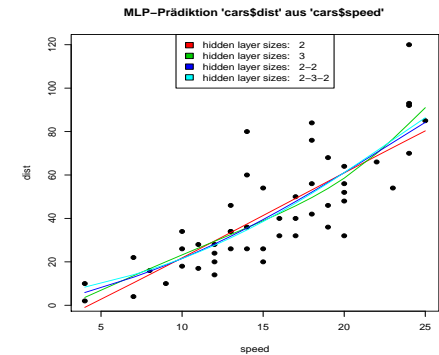
```
AANN <-function (X, k=2, sd=1/3, ...)
{
  require (neuralnet)
  X <- scale(X) * sd

  vn <- colnames(X)
  iv <- paste (vn, collapse="+")
  ov <- paste (paste (
    vn, 1, sep="."), collapse="+")
  fo <- formula (paste (iv, "~", ov))

  o <- neuralnet (formula=fo,
    data.frame(X,X), hid=k, ...)
  compute(o,X)$neurons[[2]][, -1]
}
```

## Mehrschichtenperzeptron (MLP)

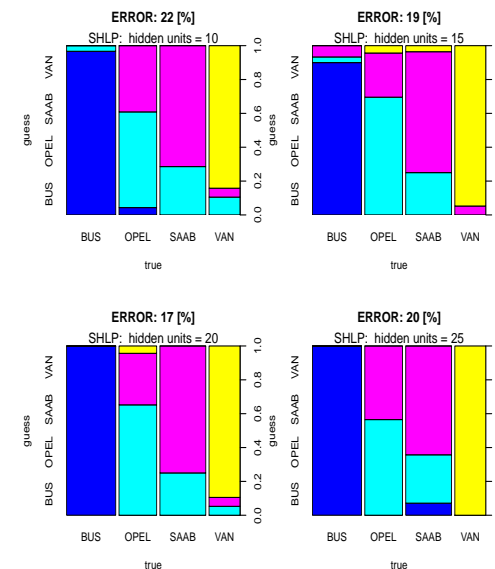
```
neuralnet::neuralnet (formula, data, hidden=1, algorithm='rprop+', ...)
```



Error: 1.621215 Steps: 23

## Dreischichtenperzeptron (SHLP)

```
nnet::nnet (x, y, size, Wts, mask, linout, entropy, softmax, censored,
  skip=FALSE, MaxNWts=1000, ...)
```



## Funktionsargumente

```
formula Output/Input-Variablen
data Datensatz
hidden Zwischenschichtengrößen
algorithm {backprop, sag, slr
           rprop+, rprop-}
return Objekt der Klasse nn
```

## R-Programmcode

```
require (neuralnet)
s <- 100
H <- list (2, 3, c(2,2), c(2,3,2))
plot (cars, pch=19)
for (j in seq_along(H)) {
  o <- neuralnet (
    dist~speed, cars/s, H[[j]])
  y <- compute(o,cars$speed/s)$net
  lines (cars$speed, y*s, col=1+j)
}
plot (o, rep="best")
```

## Funktionsargumente

```
x, y Matrix/Dataframe für E/A
formula, data (Formelinterface)
size verborgene Neuronen
skip Kurzschlusskanten?
Wts, MaxNWts Startwerte, #max
return Objekt der Klasse nnet
```

## R-Programmcode

```
require (nnet)
for (h in 5*2:5) {
  o <- nnet (CLASS~., df[[1]],
    Wts=sin(1:nwts),
    MaxNWts=nwts,
    size=h)
  guess <- predict (o, df[[2]],
    type="class")
  plot (data.frame(true,guess),
    col=3+1:4)
}
```

Modellneuronen und ihre Verschaltung

Feed-Forward-Netzwerke

Klassisches Perzeptron

Mehrschichtenperzeptron

MLPs für Prädiktion, Extraktion, Klassifikation

Rekurrente Netzwerke

## Hopfield-Netze

Vollständig verschaltetes CAM mit symmetrischer Gewichtmatrix

### Definition

Das dynamische System mit dem Zustandsübergang

$$x_k(t+1) = \text{sign} \left( \sum_{j=1}^N w_{jk} \cdot x_j(t) \right), \quad t \in \mathbb{N}, \quad k = 1, \dots, N$$

mit der symmetrischen Gewichtmatrix **W** heißt **Hopfield-Netz**. Es definiert eine (partielle) Abbildung

$$\mathbf{u} : \begin{cases} \{-1, +1\}^N & \rightarrow \{-1, +1\}^N \\ \mathbf{z} & \rightarrow \lim_{t \rightarrow \infty} \mathbf{x}(t) \mid_{\mathbf{x}(0)=\mathbf{z}} \end{cases}$$

im Raum der  $N$ -dimensionalen Binärmuster.

## CAM — Content-Addressable Memory

Assoziatives Gedächtnis: Zugriffsschlüssel  $\subseteq$  Datensatz

- **Getaktete** Neuronen mit **diskreten** Zuständen

$$x_k(t) \in \{-1, +1\}, \quad t \in \mathbb{N}, \quad k \in \mathcal{N}$$

- Nichtlineare Systemdynamik  $\mathbf{x}_t \mapsto \mathbf{x}_{t+1}$

$$x_k(t+1) = \text{sign} \left( \sum_j w_{jk} \cdot x_j(t) \right)$$

- Berechnung eines **Gleichgewichtszustandes** (Äquilibrium) **u**:

$$\mathbf{u}(\mathbf{x}_0) \stackrel{\text{def}}{=} \lim_{t \rightarrow \infty} \mathbf{x}(t)$$

- Zustandsraum  $\{0, 1\}^N$  zerfällt in **Attraktionsbassins**:

$$\mathbf{x} \cong \mathbf{y} \iff \mathbf{x} \xrightarrow{\infty} \mathbf{u} \xleftarrow{\infty} \mathbf{y}$$

## Hebb'sche Lernregel

### Lemma

Zur Speicherung von  $K$  „Prototypen“  $\mathbf{y}_1, \dots, \mathbf{y}_K$  sind die Gewichte  $w_{jk}$  proportional zum Aktivierungsratenprodukt der vernetzten Neuronen zu wählen:

$$w_{jk} = \sum_{\lambda=1}^K y_{\lambda,j} \cdot y_{\lambda,k}$$

Unter der Voraussetzung

$$K \leq 0.138 \cdot N$$

hinreichender **Kapazität** ist dann jeder Prototyp ein Äquilibrium des gelernten Netzwerks.

## Hopfield-Netze — Beispielanwendung

(John S. Denker, 1986)

## Beweis.

Für Gleichgewichtspunkte  $\mathbf{x} = \text{sign}(\mathbf{W}\mathbf{x})$  gilt:

$$\sum_j w_{jk} x_j = \sum_j \sum_{\lambda} y_{\lambda,j} y_{\lambda,k} x_j$$

Für einen Prototypen  $\mathbf{x} = \mathbf{y}_{\nu}$  gilt insbesondere

$$\sum_j w_{jk} y_{\nu,j} = \sum_j \sum_{\lambda} y_{\lambda,j} y_{\lambda,k} y_{\nu,j} = N \cdot y_{\nu,k} + \sum_j \sum_{\lambda \neq \nu} y_{\lambda,j} y_{\lambda,k} y_{\nu,j}$$

 $\mathbf{y}_{\nu}$  ist Gleichgewichtspunkt, wenn der rechte Term betragsmäßig kleiner als  $N$  ist.  $\square$ 

## Telefonbucheinträge

(25 Zeichen/Zeile, 5 Bit/Zeichen)

John Stewart Denker	8128
Lawrence David Jackel	7773
Richard Edward Howard	5952
Wayne P. Hubbard	7707
Brian W. Straughn	3126
...	...
...	...
...	...
...	...
...	...
John Henry Scofield	8109

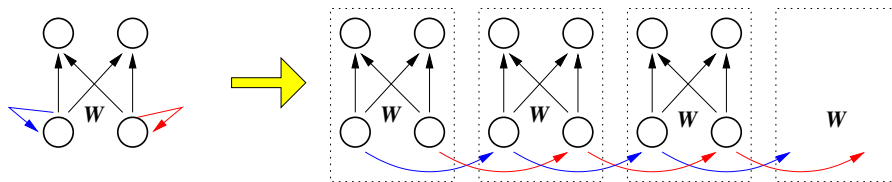
## Erfolgreicher Zugriff

t	$E = -\mathbf{x}^T \mathbf{W} \mathbf{x}$	CAM-Zustand	
0.0	-0.0	john s	
0.2	-0.0784	john sdewirubneoimv	8109
0.4	-0.8426	john sdewirtbneimv	8129
0.6	-0.8451	john sdewirtbneimv	8129
0.8	-0.8581	john sdewirt nenkmv	8128
1.0	-0.9099	john sdewart denker	8128
1.2	-0.9824	john stewart denker	8128

## Erfolgreicher Zugriff

t	$E = -\mathbf{x}^T \mathbf{W} \mathbf{x}$	CAM-Zustand	
0.0	-0.0	garbage	
0.2	-0.0244	garbagee lafj naabd	5173
0.4	-0.6280	garbaged derjd naabd	7173
0.6	-0.6904	garbaged derjd naabd	7173
0.8	-0.7345	gasbafed derjd naabd	7173
1.0	-0.7595	gasbabad derjd naabd	7173
1.2	-0.7709	fasjebad derjd naabd	7173
1.4	-0.8276	fasjebad derjd naabd	7173
1.6	-0.8282	fasjeb d derjd naabd	7173

## Zeitverzögerungsnetzwerke

TDNN — *time-delayed neural networks*

## Simulation rekurrenter Netzwerke

Das Verhalten eines rekurrenten Netzwerks mit ausschließlich **einfachen Zyklen**  $k \prec k, k \in \mathcal{N}$  kann im Rahmen eines **T-Takte-Gedächtnisses** durch ein Feed-Forward-Netz mit  $T \cdot N$  Neuronen simuliert werden.

## Zusammenfassung (8)

1. Das **Modellneuron** kombiniert seine **Eingangssignale** und berechnet daraus sein **Aktivierungsniveau**.
2. Die **Verschaltung** der Modellneuronen ist durch das Nullenmuster der **Gewichtsmatrix** synaptischer Verbindungen definiert.
3. **Rekurrente** Netzwerke sind durch die Existenz **zyklischer Verbindungen** charakterisiert und beschreiben **dynamische Systeme**, deren Verhalten durch die **Attraktionsgebiete** ihrer Fixpunktzustände definiert ist.
4. **Nichtrekurrente** Netzwerke beschreiben pure Vorwärtsberechnungen und definieren durch die Zustände ihrer **Eingabe-** und **Ausgabeneuronen** eine **nichtlineare Vektorabbildung**.
5. Das **einstufige Perzeptron** ist nicht in der Lage, die **logische XOR-Funktion** zu lernen.
6. Ein **mehrstufiges Perzeptron** (MLP) mit mindestens vier Neuronenschichten gilt als (asymptotisch) **universeller Approximator**.
7. Jedes MLP mit **differenzierbarer** Aktivierungsfunktion kann mittels **Error-Backpropagation**-Algorithmus seine Gewichtsmatrix im Sinne **kleinster Fehlerquadrate** auf eine Lerndatenprobe optimieren.
8. Für die **numerische Klassifikation** verwendet man MLPs mit  $D$  Eingabeneuronen,  $K$  Ausgabeneuronen und einer der anvisierten **Modellkapazität** angemessenen Anzahl **verborgener Neuronen**.