

Neural Nets, Background Back-Propagation

Neuronale Netze u. Anwendungen, 17. Jan. 2012



Overview

- Backpropagation
 - Implementation
 - Speed, Training
 - Generalization
 - Choice of Architecture
 - Algorithms



Training Time

- Backpropagation is Computer Intensive during **Training**
very efficient during evaluation/testing
- How to Improve Training Time
 - Parallel Hardware (?)
 - Efficient Implementation
 - Faster Gradient Descent Search
 - Selective Choice of Patterns
 - Efficient Architectures



Back-Propagation of Error

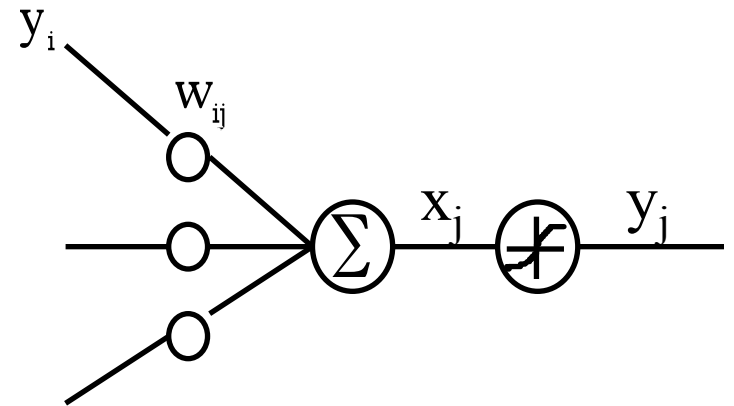
$$E = \frac{1}{2} \sum_j (y_j - d_j)^2 \quad y_j = \frac{1}{1 + e^{-x_j}} \quad x_j = \sum_i y_i w_{ij}$$

$$1). \quad \frac{\partial E}{\partial y_j} = y_j - d_j$$

$$2). \quad \frac{\partial E}{\partial x_j} = \frac{\partial E}{\partial y_j} \frac{\partial y_j}{\partial x_j} = \frac{\partial E}{\partial y_j} y_j [1 - y_j]$$

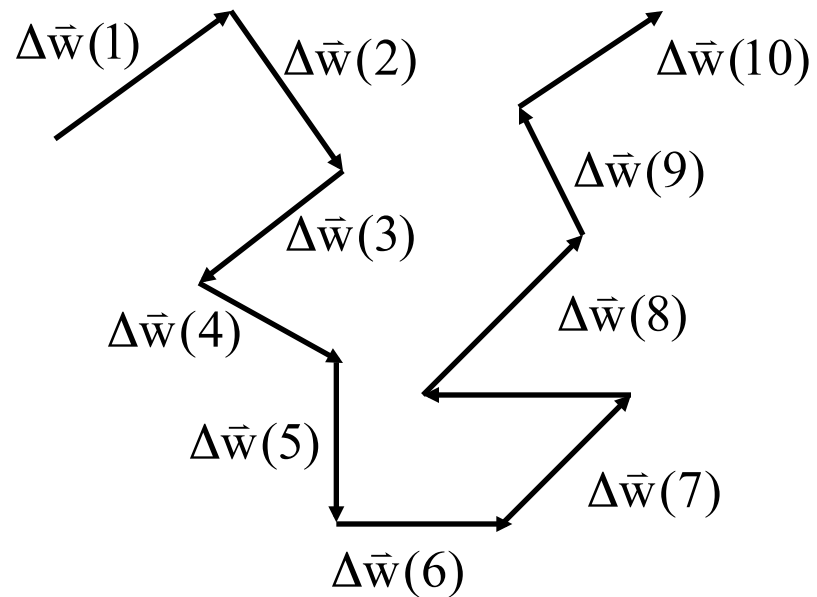
$$3). \quad \frac{\partial E}{\partial w_{ij}} = \frac{\partial E}{\partial x_j} \frac{\partial x_j}{\partial w_{ij}} = \frac{\partial E}{\partial x_j} y_i$$

$$4). \quad \frac{\partial E}{\partial y_i} = \sum_j \frac{\partial E}{\partial x_j} \frac{\partial x_j}{\partial y_i} = \sum_j \frac{\partial E}{\partial x_j} w_{ij}$$



Momentum

$$\Delta w_{ij}(t) \sim \frac{\partial E}{\partial w_{ij}}$$



$$\Delta w_{ij}^*(t+1) = \Delta w_{ij}(t+1) + \alpha \Delta w_{ij}(t)$$

↑
Momentum
0...1



Überspringen von Trainingsbeispielen

- Kein BP für Trainingsbeispiele, die bereits gut gelernt sind

Vorwärts-Durchlauf des Beispiels X durch das Netzwerk

Fehler am Ausgang $>$ Schwellenwert?

Ja: Back-Propagation

Nein: nächstes Trainingsbeispiel



Trainingsprozedur

A	B	C	D	E	F
A_1	B_1	C_1	D_1	E_1	F_1
A_2	B_2	C_2	D_2	E_2	F_2
A_3	B_3	C_3	D_3	E_3	F_3
⋮					
A_n	B_n	C_n	D_n	E_n	F_n

$(A_1, B_1, \dots, E_n, F_n) \rightarrow$ 1 weight update

└───┬───┘
└───┘ langsam aber sicher!



Trainingsprozedur

$A_1 \longrightarrow 1$ weight update

$B_1 \longrightarrow 1$ “ “

.....

$E_n \longrightarrow 1$ “ “

$F_n \longrightarrow 1$ “ “

\longrightarrow schnell, aber riskant!

$(A_1 \dots F_1) \longrightarrow 1$ weight update

$(A_2 \dots F_2) \longrightarrow 1$ weight update



Dynamische Anpassung der Lernrate ε

(Franzini)



$$\cos \Theta = \frac{\sum_{i,j} (\Delta w_{ij}(t-1) \cdot \Delta w_{ij}(t))}{\sqrt{\sum_{i,j} (\Delta w_{ij}(t-1))^2} \cdot \sqrt{\sum_{i,j} (\Delta w_{ij}(t))^2}}$$

$$\cos \Theta = 1 \quad (\Theta = 0^\circ, 360^\circ, \dots)$$

$\Rightarrow \varepsilon$ größer machen

$$\cos \Theta = -1 \quad (\Theta = 180^\circ)$$

$\Rightarrow \varepsilon$ kleiner machen

$$\varepsilon(t) = \varepsilon(t-1) \cdot \text{const} \cdot \frac{\cos \Theta + 1}{2}$$

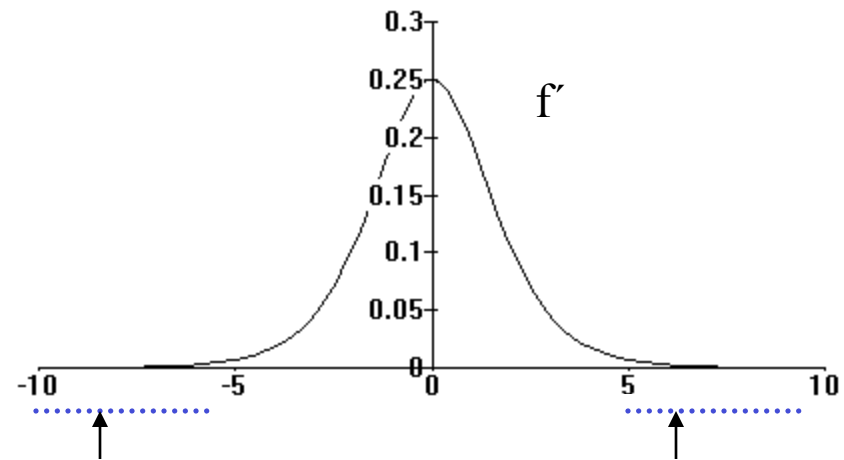
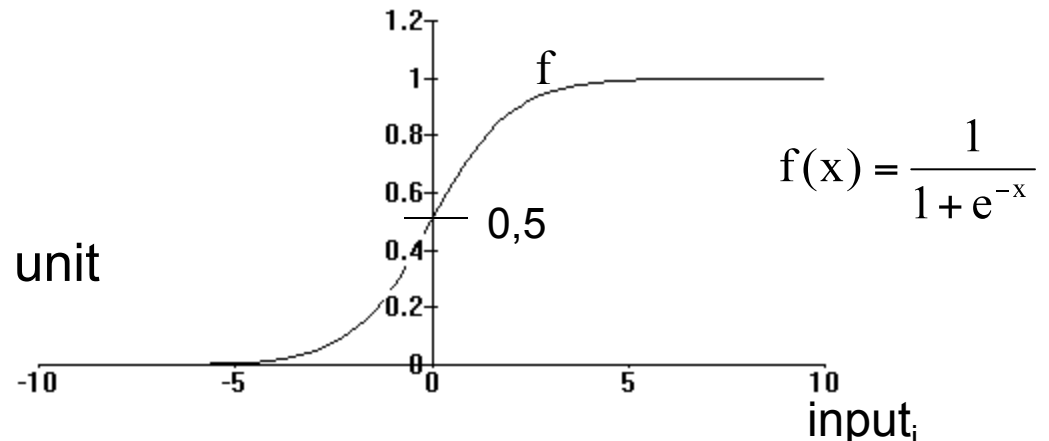


Beschleunigung der Lernphase

$$\text{BP: } \Delta w_{ij} = \varepsilon \cdot \delta_j \cdot O_i$$

↑
output i-tes unit

$$\delta_j \sim f'(\text{input unit}_j)$$



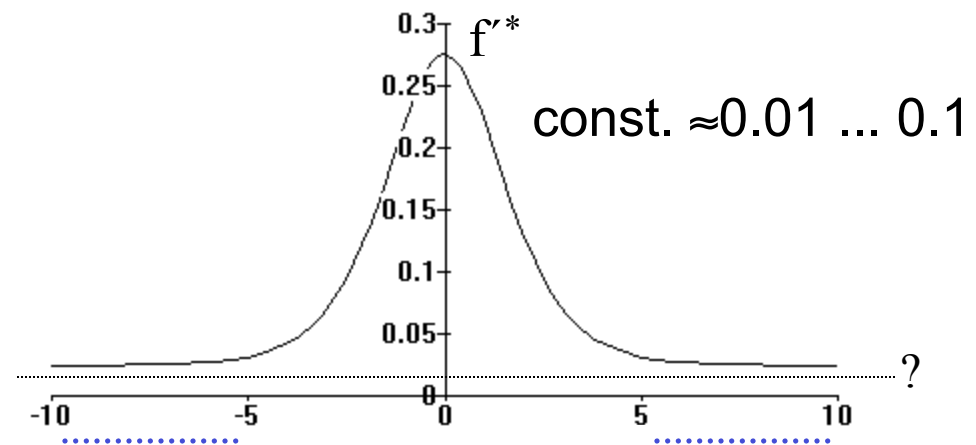
f' sehr klein
 $\Rightarrow \Delta w_{ij}$ sehr klein

Beschleunigung der Lernphase

Lösung:

(S. Fahlmann)

$$f'^* = f' + \text{const.}$$



- i.A. schnelle Konvergenz
- besonders wirksam bei hochgradig nichtlinearen Klassifizierungsproblemen (XOR)
- manchmal schlechtere Generalisierung auf den Test-Daten

Der Quickprop Algorithmus

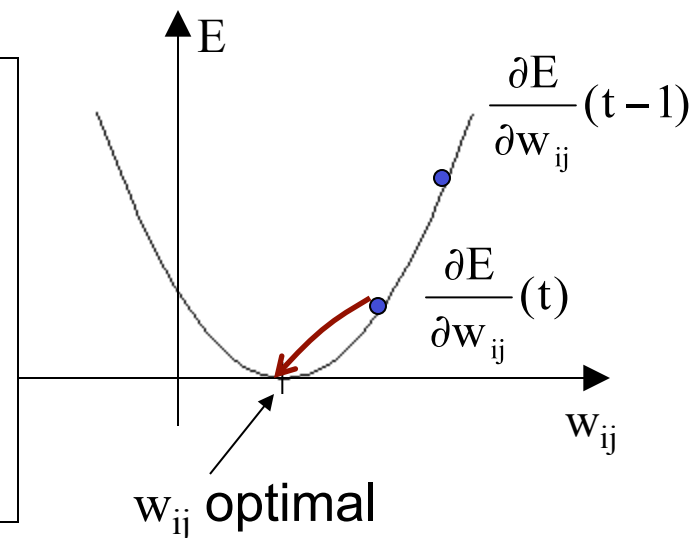
Ziel: schnelles Training eines MLP

Zwei „riskante“ Annahmen:

- 1.) die Fehler/weight Kurve hat die Form einer Parabel
- 2.) die Parabeln für jedes weight w_{ij} sind unabhängig voneinander

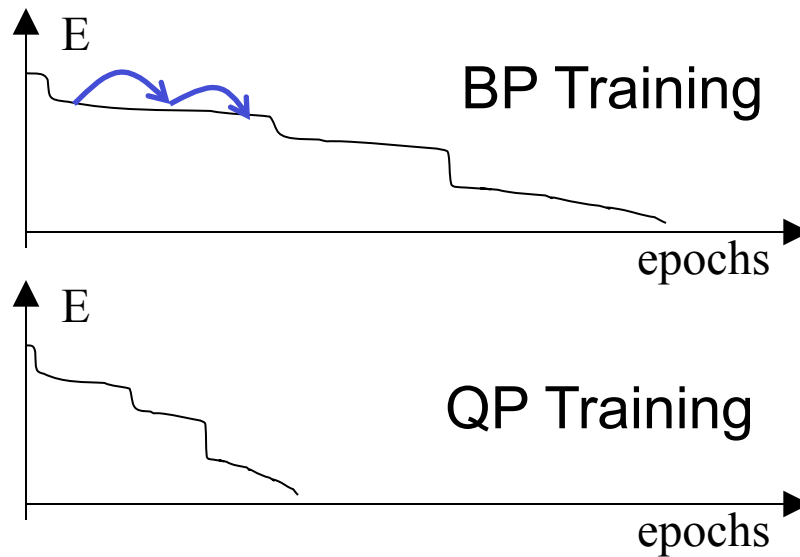
$$\Delta w(t) = \frac{s(t)}{s(t-1) - s(t)} \cdot \Delta w(t-1)$$

mit $s(t) = \frac{\partial E}{\partial w_{ij}}(t)$



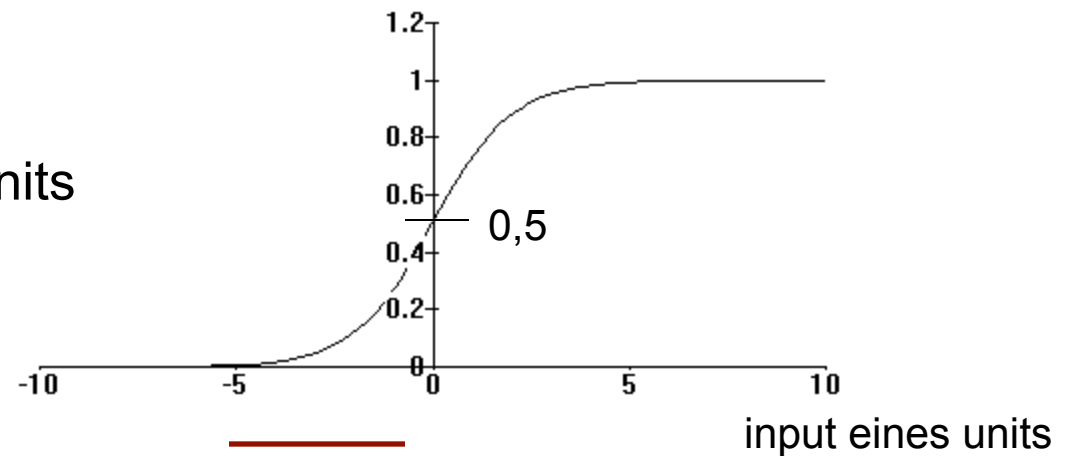
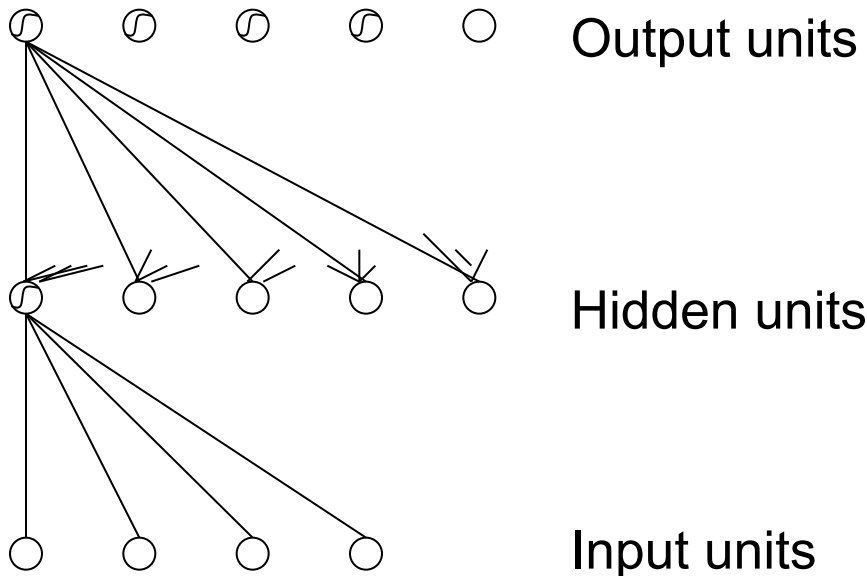
Quickprop

- Annahme 2 ist nicht richtig,
aber: durch das iterative Einstellen der weights
werden die gemachten Fehler aufgehoben
- funktioniert sehr gut für hochgradig nichtlineare Probleme
wie z.B. XOR, die sich so verhalten:



Initialisierung mit Zufallszahlen

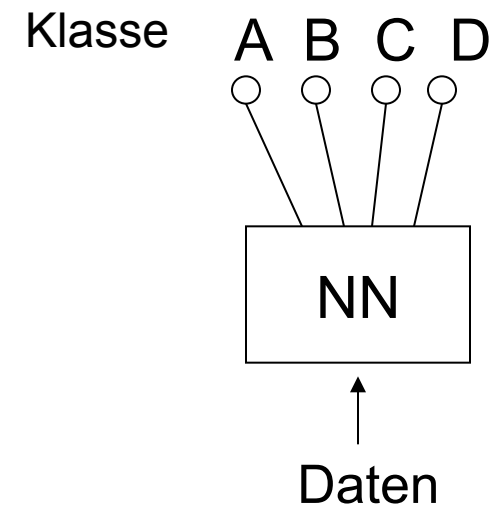
-0.3 ... +0.3



Klassifikatoren: Ausgangsrepräsentation

Einfachst mögliche Ausgangsrepräsentation:

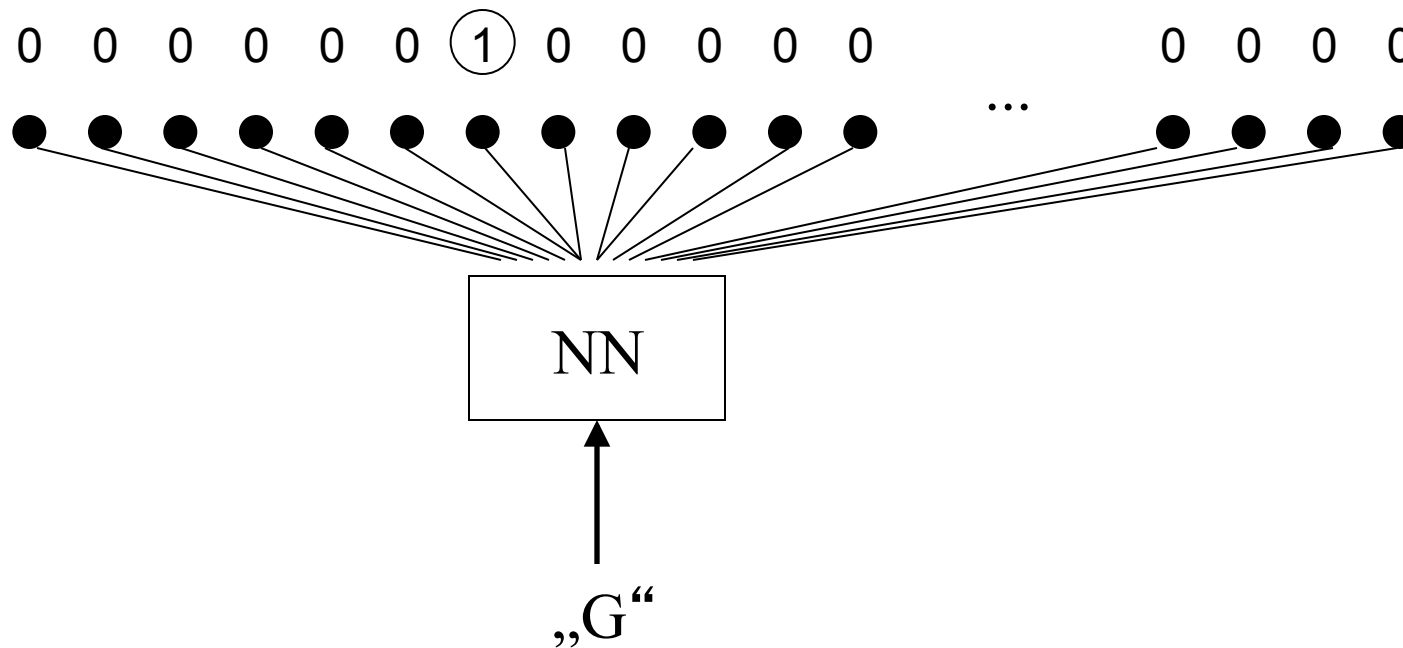
Klasse A	1	0	0	0
Klasse B	0	1	0	0
Klasse C	0	0	1	0
Klasse D	0	0	0	1



- Vorteile:
- die wahrscheinlichste Klasse ist leicht zu finden (Maximum suchen)
 - Abschätzung der Verlässlichkeit der Klassifikation ist möglich

Probleme bei vielen Klassen

Angenommen: Alphabet-Erkenner
=> 26 Klassen



am einfachsten, wenn alle output units Null gesetzt werden!

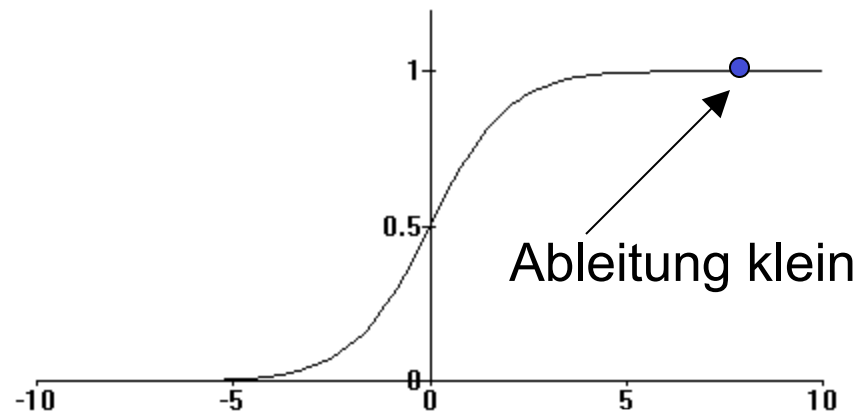
Probleme mit dem MSE

$$E = \text{MSE} = \sum_j (y_j - \text{target}_j)^2$$

angenommen: $y_j = 0,99$
 $\text{target}_j = 0,0$

=> großer Beitrag zu E

aber:

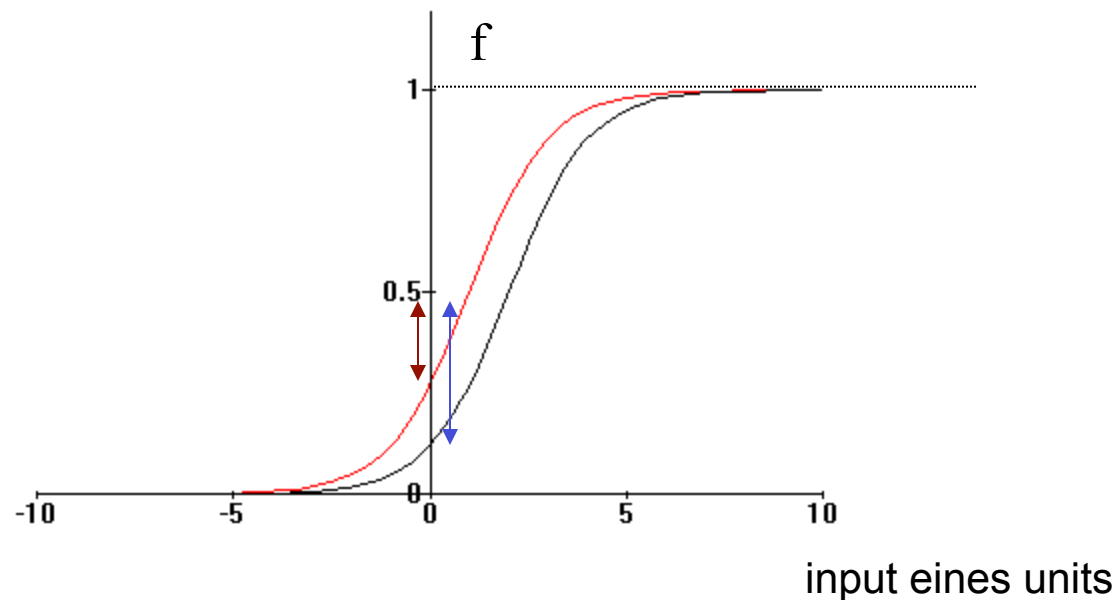


=> Δw klein

=> langsameres Lernen, obwohl E groß ist



Asymmetrische Sigmoid Funktion



- stetig
- beschränkt
- Ableitung sollte einfach auszurechnen sein
- monoton

Andere Fehlerfunktionen

$$E = - \sum_j \ln(1 - (y_j - \text{target}_j)^2) \quad (\text{z.B. Mc Clelland})$$

$$\frac{\partial E}{\partial y_j} = - \frac{1}{1 - (\text{target}_j - y_j)} + \frac{1}{1 + (\text{target}_j - y_j)}$$

Multiplikation mit $f' = y_j(1 - y_j)$:

$$\frac{\partial E}{\partial y_j} = \frac{(1 - y_j)^2}{y_j - 2} \quad \text{Falls } \text{target}_j = 1 \quad \frac{\partial E}{\partial y_j} = \frac{y_j^2}{y_j + 1} \quad \text{Falls } \text{target}_j = 0$$

$$\Rightarrow \frac{\partial E}{\partial y_j} \quad \text{maximal, wenn } y_j - \text{target}_j = 1$$



Weitere Fehlerfunktionen

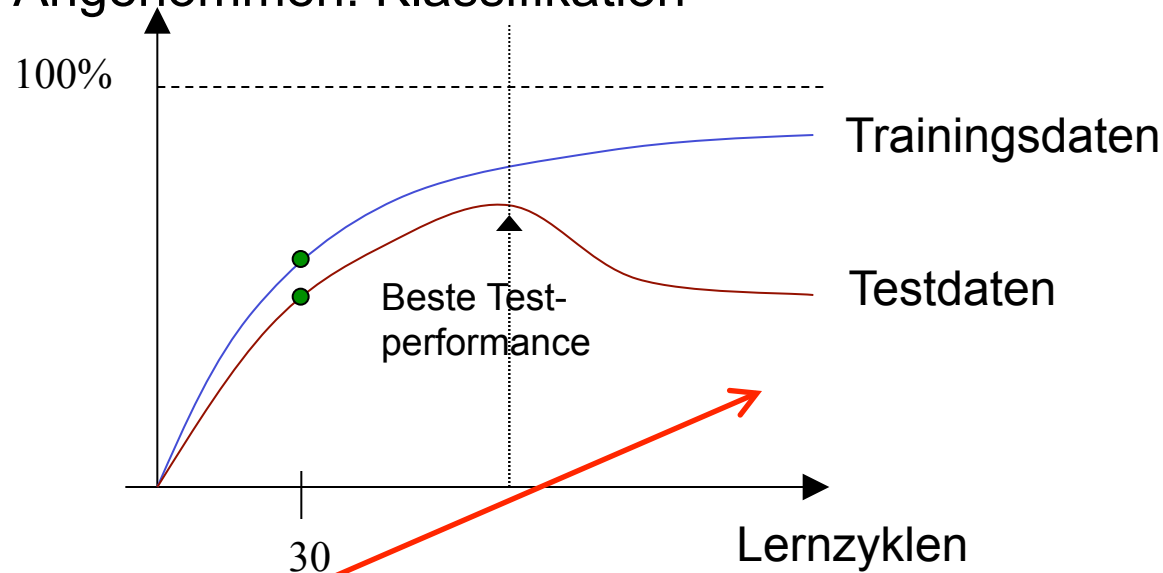
- Backpropagate on Error from Best and Correct output



Generalisierungsfähigkeit:

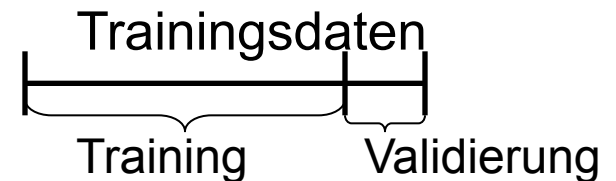
Fähigkeit, das aus den Trainingsdaten Erlernte auf neue Daten (Testdaten) anzuwenden.

Angenommen: Klassifikation



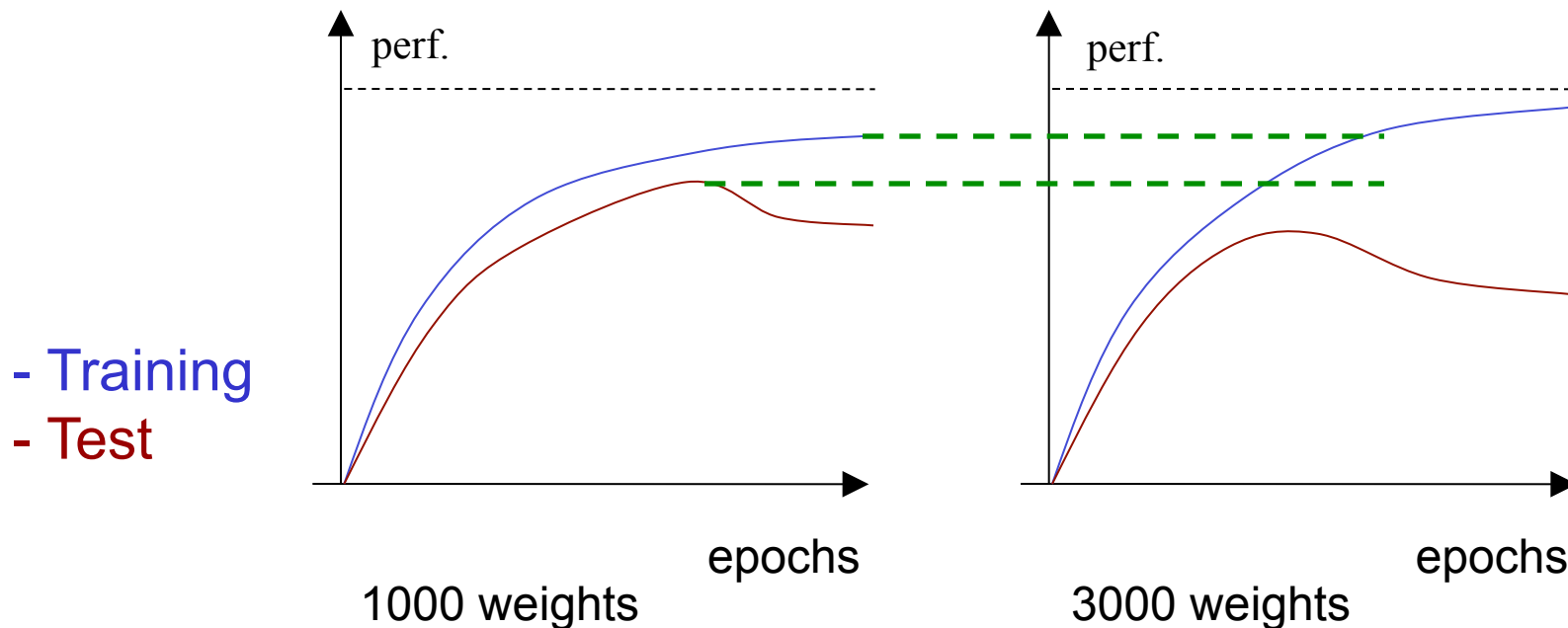
**Auswendiglernen
„Overfitting“**

→ Eine einfache Lösung: Validation Data



Zwei Gründe für schlechte Generalisierung?

- 1.) Overfitting/Overtraining (zu lange trainiert)
- 2.) Zu viele Parameter (weights) bzw. zu wenig Trainingsmaterial
- 3.) falsche Netzwerkstruktur (neu)



Generalisierung in linearen Systemen

Für lineare Systeme gilt:

(Akaike, 1970: Statistical Predictor Identification
Ann. Inst. Stat. Math., 22:203 oder Barron oder Eubank)

$$\langle \varepsilon_{\text{test}} \rangle_{\xi \zeta'} = \langle \varepsilon_{\text{train}} \rangle_{\xi} + 2\sigma^2 \frac{p}{n}$$

ξ : Training Set
 ζ' : Test Set
 p : # Parameter
 n : # Trainingsbeispiele

Effektive Varianz
des Rauschens



Generalisierung in nichtlinearen Systemen

(Moody, 1992)

$$\begin{aligned} \langle \varepsilon_{\text{test}} \rangle_{\zeta\zeta'} &= \langle \varepsilon_{\text{train}}(\lambda) \rangle_{\xi\xi'} \\ &\approx \langle \varepsilon_{\text{train}}(\lambda) \rangle_{\xi} + 2\sigma_{\text{eff}}^2 \frac{p_{\text{eff}}(\lambda)}{n} \end{aligned}$$

σ_{eff}^2 : effektive Varianz des Rauschens

p_{eff} : effektive Anzahl der Parameter

$$p_{\text{eff}} \leq p$$

λ : Regulierungsparameter



Methoden zur Verbesserung der Generalisierung

- 1.) Reduzierung der Komplexität des Netzwerkes durch Regularisierung
 - Weight Decay
 - Weight Elimination
 - Optimal Brain Damage
 - Optimal Brain Surgeon
 -
- * Ein großes Netzwerk wird daran gehindert, die Daten auswendig zu lernen.



Methoden zur Verbesserung der Generalisierung

2.) Schrittweises Vergrößern eines zu kleinen Netzes (konstruktiv)

- Cascade Correlation
- Meiosis Netzwerke
- ASO (Automativ Structure Optimization)

⋮

- * Ein zu kleines Netzwerk wird solange erweitert,
bis es die Daten lernen kann.



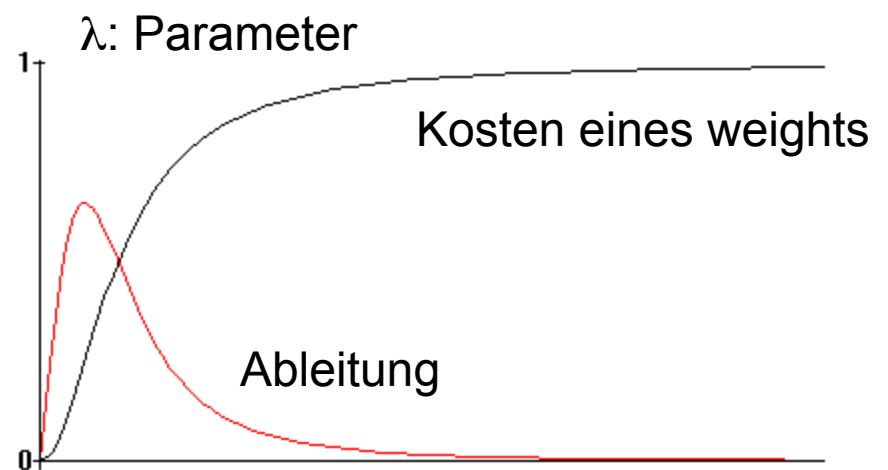
Regularisierung

- soll das „Overfitting“ verhindern

z.B. Weight-Elimination (Rumelhart, Chauvin, Denker, ...)

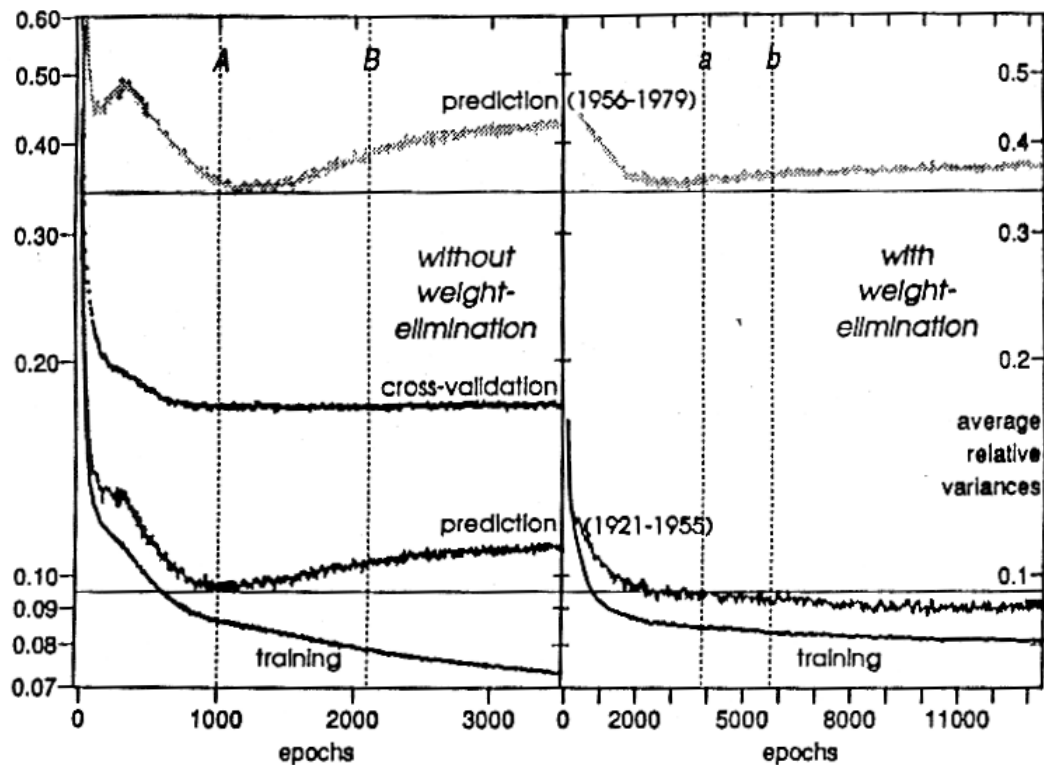
$$E = \text{MSE} + \lambda \sum_{i,j} \frac{w_{ij}^2}{1 + w_{ij}^2}$$

„Kosten“ eines weights



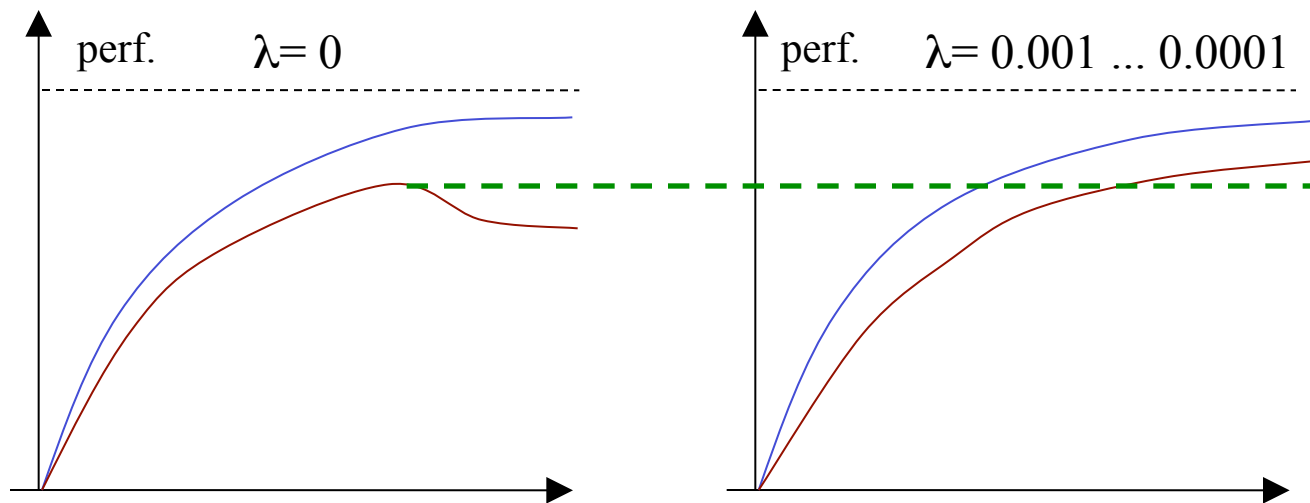
Interpretation : das Netzwerk wird für große $|w|$ bestraft!

Lernkurven eines 12-8-1 Netzwerks



The average relative single-step prediction variance are given for the training sets (solid lines) and the early (low gray lines) and the late prediction set (upper gray lines) (as well as for the cross-validation set for the network trained without weight-elimination on the left side). The vertical lines A, B, a, b) indicate different stopping points

Weight-Elimination : Auswirkungen



- langsames Lernen
- schlechtere Performanz auf den Trainingsdaten
- bessere Generalisierung

Einfachste Form:

$$\text{Error} = \text{MSE} + \lambda \sum_{i,j} w_{ij}^2$$

„Weight Decay“



Optimal Brain Damage

Idee: bestimmte Verbindungen werden aus dem Netzwerk entfernt um die Komplexität des Netzes zu beschränken und so Overfitting zu vermeiden.

Einfachste Form:

Es werden diejenigen Verbindungen entfernt, die nach dem Training sehr kleine Beträge $|w_{ij}|$ haben.

Richtig:

Es werden diejenigen Verbindungen entfernt, die dabei den geringsten Einfluß auf den Fehler (MSE oder ä.) haben, d.h. die unwichtigsten Verbindungen.

dazu: $\frac{\partial E}{\partial w_{ij}^2}$ berechnen
-> zeitaufwendig (schwierig)

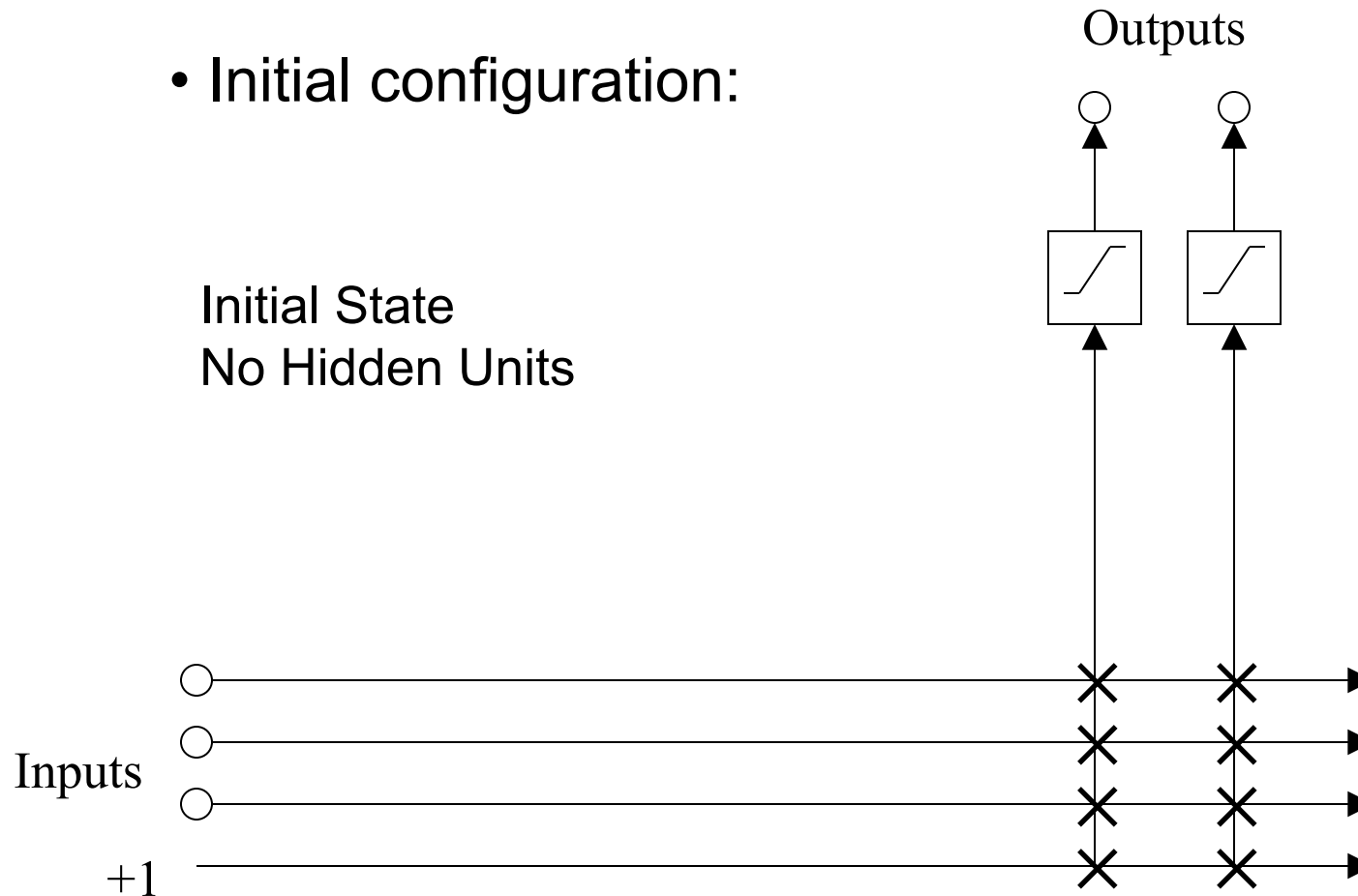


Cascade Correlation

(Fahlman, Lebiere)

- Initial configuration:

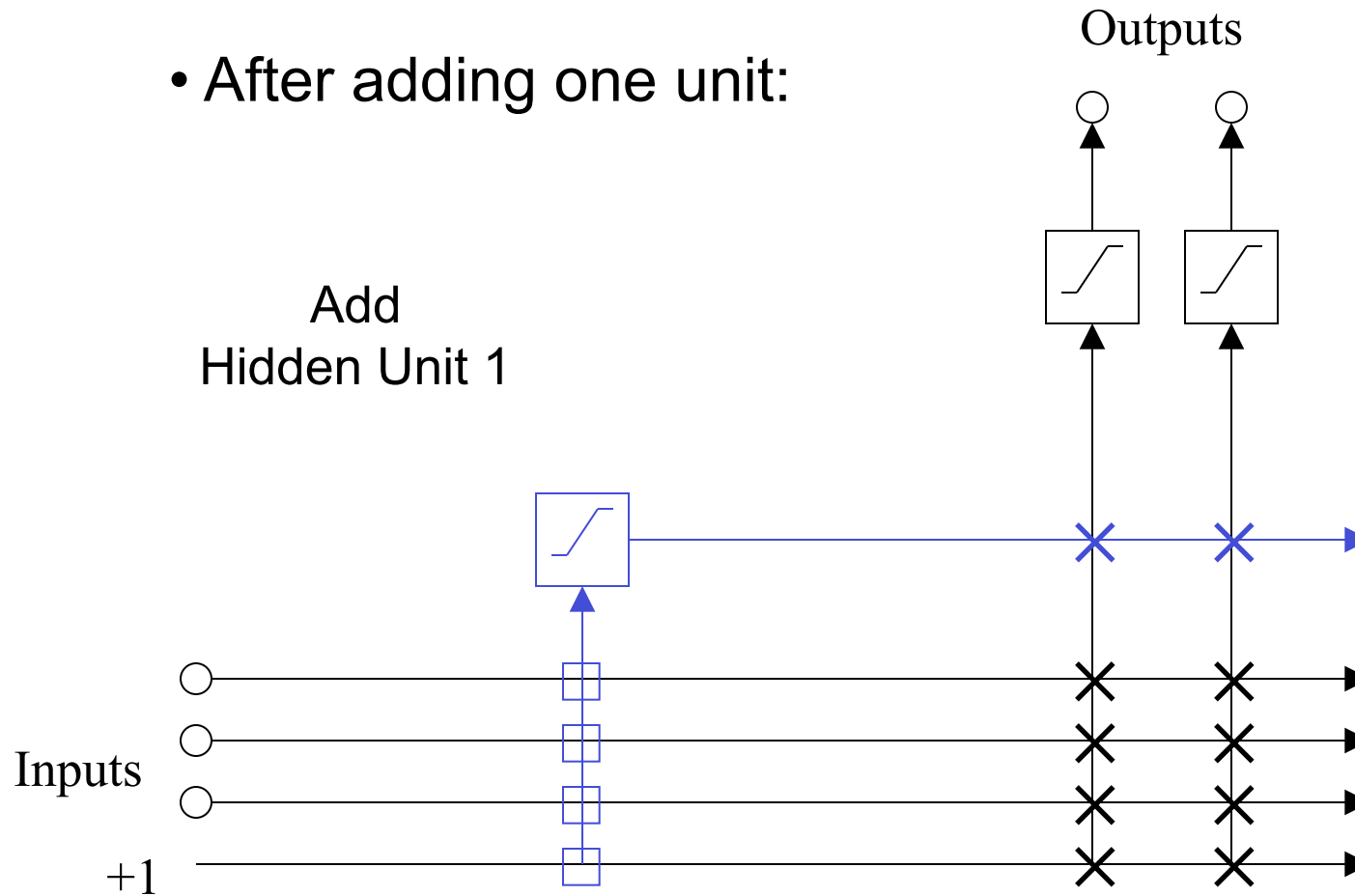
Initial State
No Hidden Units



Cascade Correlation

(Fahlman, Lebiere)

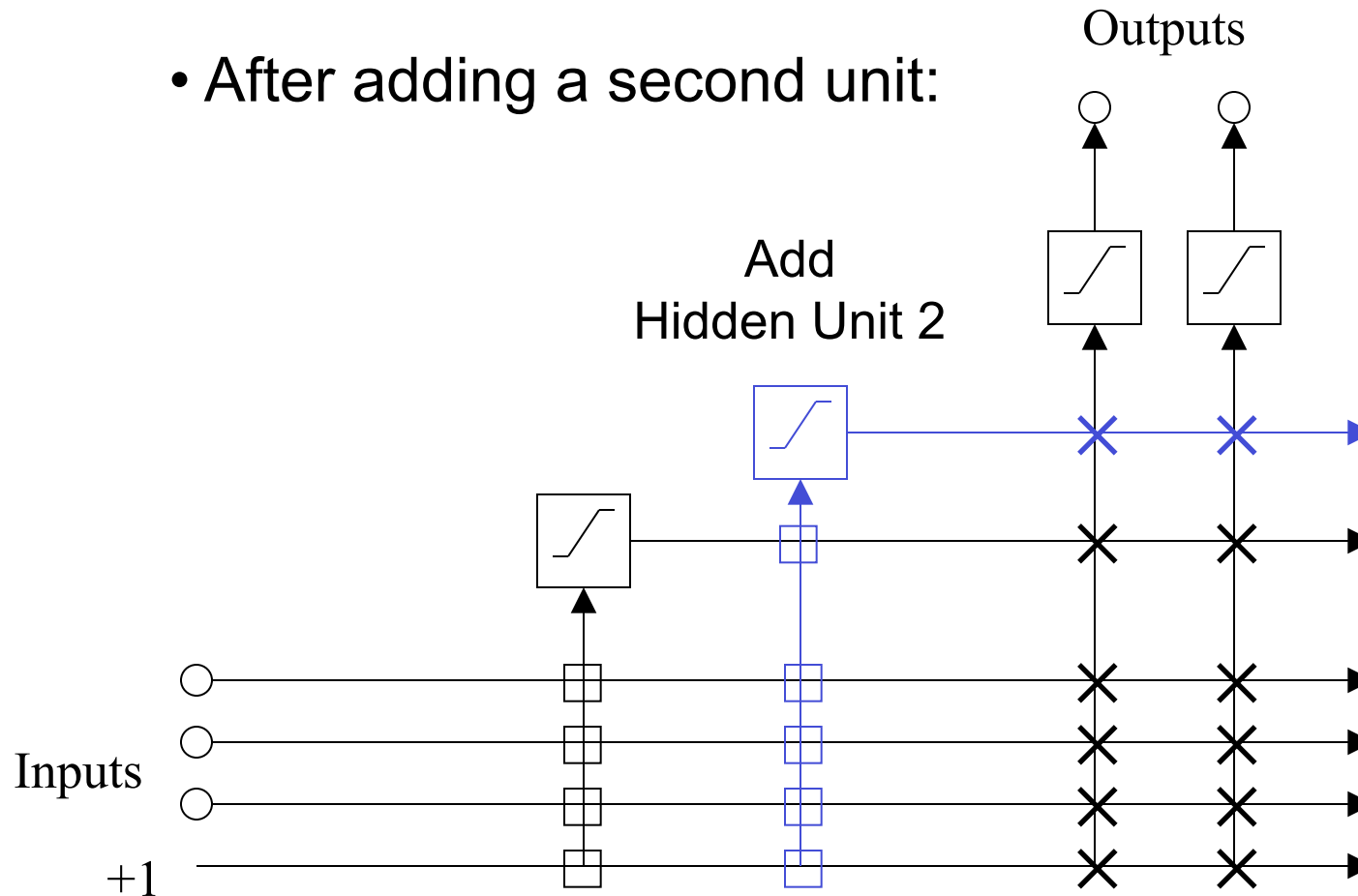
- After adding one unit:



Cascade Correlation

(Fahlman, Lebiere)

- After adding a second unit:



Cascade Correlation

- Hinzufügen eines hidden units:
 - input connections von allen input units und allen bisherigen hidden units
 - Zunächst werden nur diese Verbindungen adaptiert.

Ziel: Maximieren der Funktion

$$S = \sum_o \sum_p (V_p - \bar{V})(E_{p,o} - \bar{E}_o)$$

p: Trainingsmuster
o: Index der output units

$$\frac{\partial S}{\partial w_i} = \sum_{p,o} \sigma_o(E_{p,o} - \bar{E}_o) f'_p I_{i,p}$$

Vorzeichen der Correlation
zwischen der Aktivierung
des candidate units und output o

Ableitung der
sigmoid Funktion

Input des
candidate units

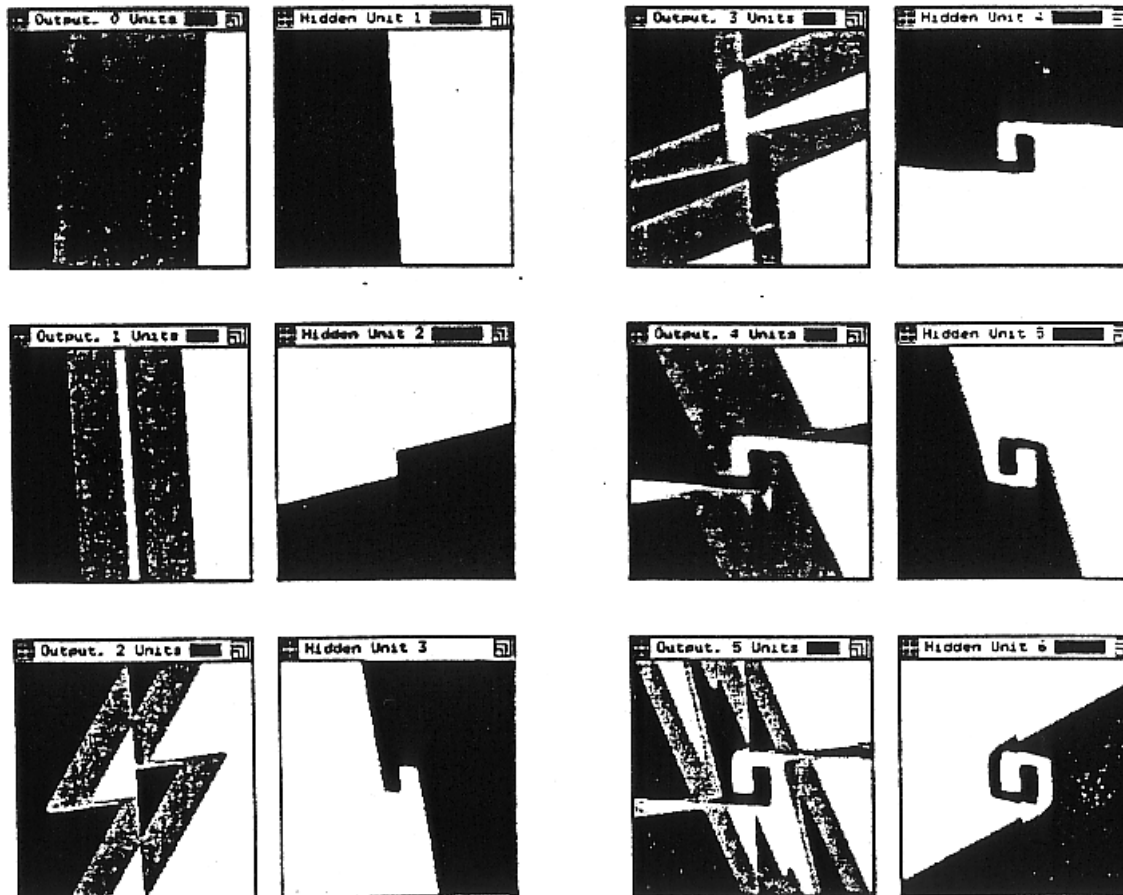


Zusammenfassung Cascade Correlation

- Nicht notwendig, die Anzahl Hidden Units oder Hidden Layers empirisch zu bestimmen.
- Kann tiefe Netze erzeugen ohne dramatischen Slowdown.
- Zu jedem Zeitpunkt wird nur eine Ebene von Verbindungen trainiert.
- Lernt sehr schnell.
- Inkrementelles Training.

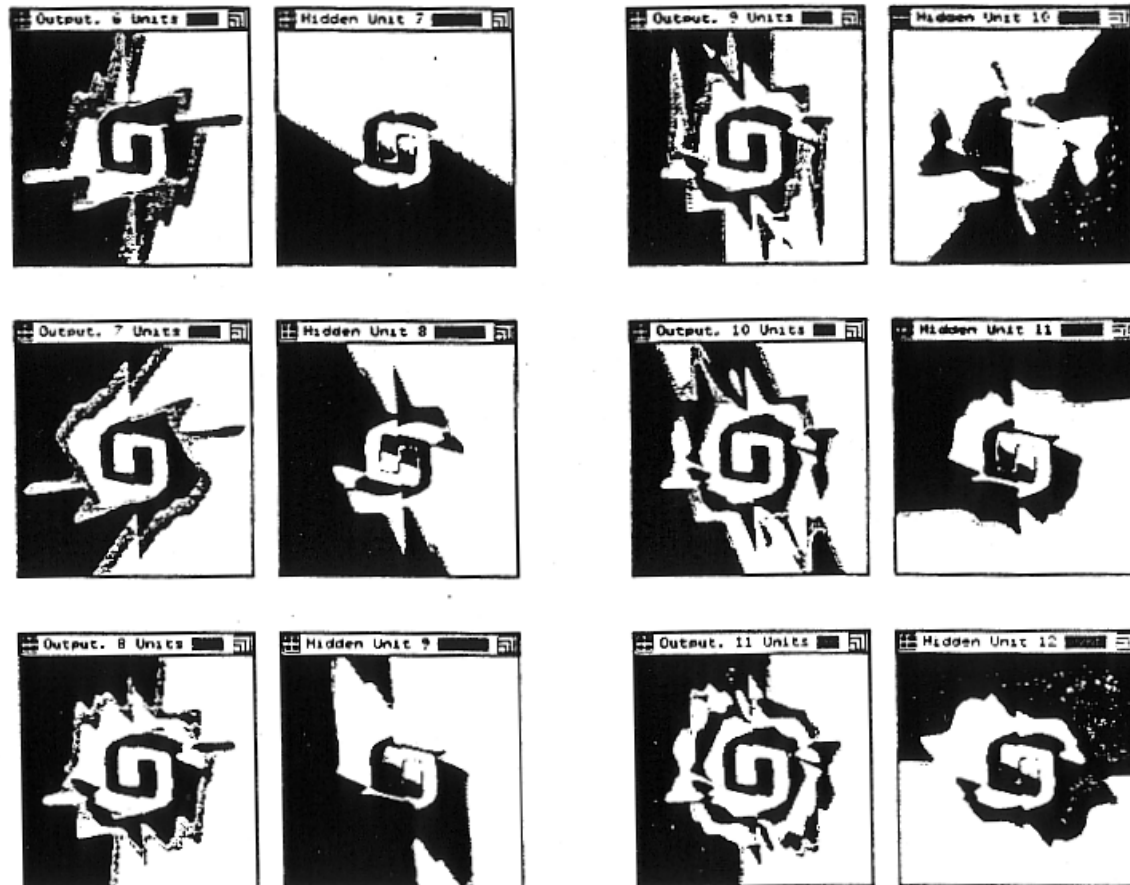


Entwicklung einer 12-hidden-unit Loesung fuer das zwei-Spiral-Problem



black: strong negative
white: strong positive

Entwicklung einer 12-hidden-unit Loesung fuer das zwei-Spiral-Problem(kont.)



black: strong negative
white: strong positive

Meiosis Netzwerke

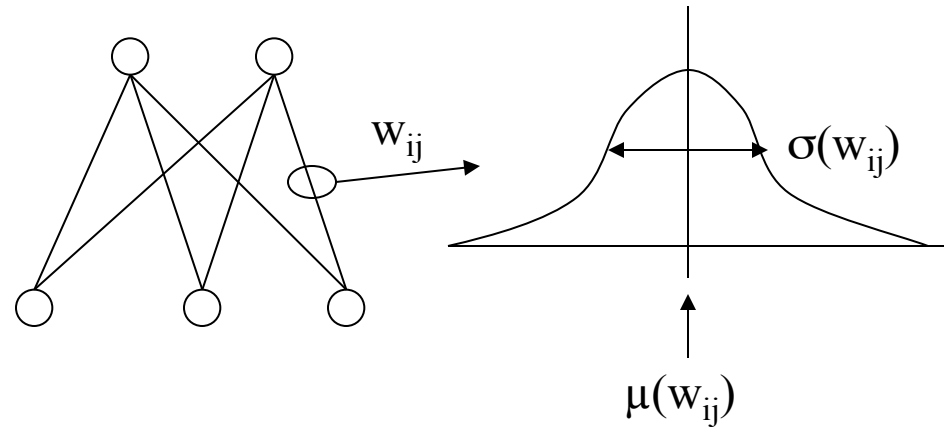
(Hanson, 1990)

Idee: Hinzufügen von hidden units abhängig von der „Unsicherheit“ der weights

$$x_i = \sum_j w_{ij}^* \cdot y_i$$

$$w_{ij}^* = \mu(w_{ij}) + \underbrace{\sigma(w_{ij})\phi(0,1)}_{\text{Gauss Noise}}$$

Gauss Noise
mit Mittelwert 0 und Varianz 1



- der Mittelwert eines Gewichts $\mu(w_{ij})$ und die Varianz $\sigma(w_{ij})$ werden gelernt

Meiosis, Teil 2

- Start mit einem hidden unit
- ein hidden unit wird aufgespalten, falls

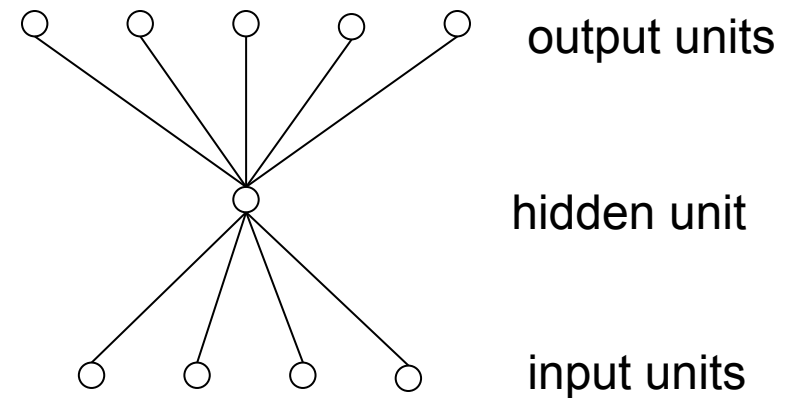
$$\frac{\sum_i \sigma_{ij}}{\sum_i \mu_{ij}} > 1.0$$

Verbindung zu
den Input Units

und

$$\frac{\sum_k \sigma_{ik}}{\sum_k \mu_{ik}} > 1.0$$

Verbindung zu
den Output Units



Probleme:

- 1.) Zusätzliche Parameter der Verbindungen (σ)
- 2.) Rausch-Generator
- 3.) Instabilitäten beim Hinzufügen von neuen hidden units