

## VIII. LERNEN UND OPTIMIEREN

- In den meisten Lernverfahren für neuronale Netzwerke ist "lernen" gleichbedeutend mit "optimieren":

Lernen = Prozess, welcher die Effizienz des Netzwerks (bezüglich einer vorgegebenen Aufgabe) optimiert.

- Wie die Effizienz eines Netzwerks gemessen oder bewertet wird, hängt von der Problemstellung ab.

Ist z.B. der gewünschte Output  $\{D_i^*\}$  für einen Satz von Lernbeispielen vorgegeben ( $r=1, \dots, N$ ), so wählt man oft den quadratischen Outputfehler als Mass für die Effizienz:

$$F = \sum_{r=1}^N F^r \quad , \quad F^r = \frac{1}{2} \sum_i (D_i^r - O_i^r)^2 \quad .$$

Da  $O_i^r = O_i^r(\underline{W})$ , ergibt sich ein Optimierungskriterium der Form

$$F(\underline{W}) = \min \quad ,$$

d.h. die Optimierung bezieht sich auf die Wahl der Verbindungsgewichte  $\underline{W} = \{W_{ij}\}$ .

- Die zwei Hauptfaktoren, die das Lernverhalten eines neuronalen Netzwerks bestimmen, sind
  - 1) die Topographie der "Fehlerlandschaft" über dem  $\underline{W}$ -Raum
  - 2) die gewählte Lernstrategie (bzw. Optimierungsstrategie).

## VIII.1 Gradientenverfahren

Die meisten neuronalen Lernverfahren (siehe z.B. "Error Backpropagation", Kapitel VIII) beruhen auf einer Gradientenmethode zur Minimierung der Fehlerfunktion  $F(W)$ :

$$W_{ij} \rightarrow W_{ij} + \Delta W_{ij}$$

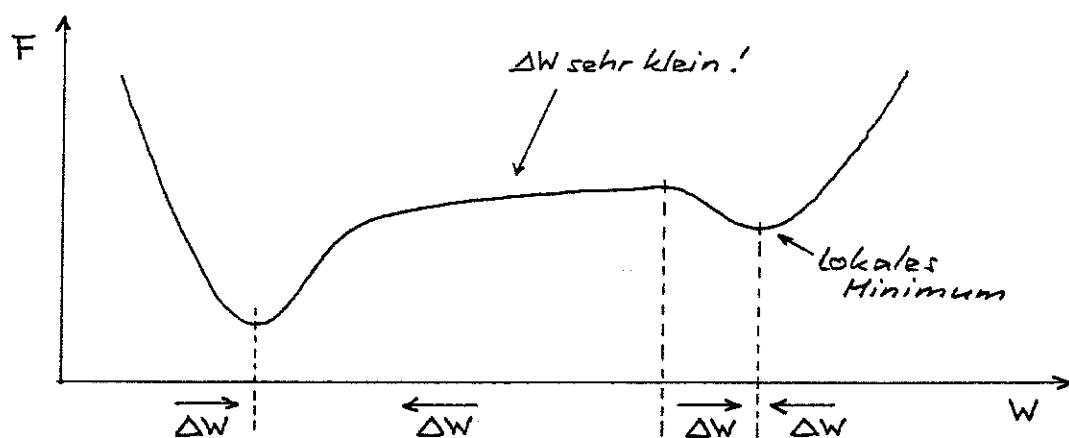
"Example-by-example"-Version:  $\Delta W_{ij} = -\eta \frac{\partial F^*}{\partial W_{ij}}$ ,

wobei die Lernbeispiele  $n$  in zufälliger Reihenfolge immer wieder präsentiert werden.

"Batch"-Version:  $\Delta W_{ij} = -\eta \frac{\partial F}{\partial W_{ij}}$ ,

d.h. die Fehler werden über alle Lernbeispiele summiert, bevor die Gewichte geändert werden.

- Probleme mit Gradientenverfahren bei komplizierten Fehlerlandschaften:
  - Langsame Konvergenz (flache Täler u. Plateaus)
  - Steckenbleiben in schlechten lokalen Minima

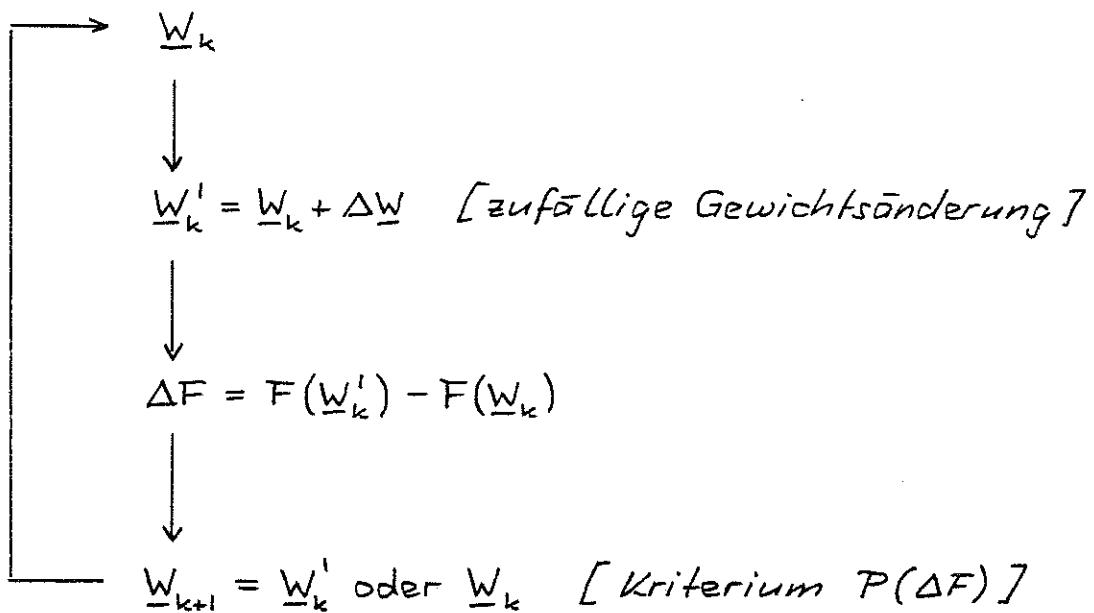


- Gradientenverfahren sind zudem nur anwendbar, wenn die Gewichte kontinuierlich variieren dürfen (d.h.  $-\infty < W_{ij} < +\infty$ ) und wenn die Aktivierungsfunktion  $f(\sum_j W_{ij} S_j)$  differenzierbar ist.

## VIII.2 Stochastische Lernstrategien

Problem:  $F(\underline{w}) = \sum_{i=1}^n F_i(\underline{w}) = \min$

- Stochastische Lernverfahren sind Strategien, die eine zufällige Suche im  $\underline{W}$ -Raum gegen ein gutes Minimum der Fehlerfunktion  $F(\underline{W})$  steuern:



- Solche Verfahren sind auch dann anwendbar, wenn die Gewichte auf diskrete Werte beschränkt sind (z.B.  $W_{ij} = \pm 1$ ) [ $\rightarrow$  Kombinatorische Optimierung], oder wenn die Aktivierungsfunktion  $f$  nicht differenzierbar ist.
  - Implementierungsoptionen:
    - "Batch" oder "Example-by-Example"
    - Wahl der Wahrsch.keitsverteilung  $\mu(\Delta \underline{w})$  für die zufälligen Gewichtsänderungen
    - Kriterium  $P(\Delta F)$  für das Akzeptieren des zufällig erzeugten Gewichtsvektors  $\underline{W}'_k$

### VIII.3 "Iterative Improvement"

→ Kriterium  $P(\Delta F)$ :

$$\underline{w}_{k+1} = \begin{cases} \underline{w}_k' & \text{falls } \Delta F \leq 0 \\ \underline{w}_k & \quad \quad \quad \Delta F > 0 \end{cases}$$

d.h. die Gewichtsänderung wird nur akzeptiert, wenn die Fehlerfunktion abnimmt.

→ Stochastisches Äquivalent zur Gradientenmethode!

→ Ähnliche Probleme mit lokalen Minima, falls  $|\Delta \underline{w}|$  beschränkt ist.

Aber Satz: Falls 1) die Suche auf einen kompakten Unterraum  $S$  des  $\underline{w}$ -Raums beschränkt wird,  
2) die Reichweite von  $\mu(\Delta \underline{w})$  größer als der Durchmesser von  $S$  ist,

dann:

$$\lim_{k \rightarrow \infty} \text{Prob}\{F(\underline{w}_k) < F_{\min}^{(S)} + \varepsilon\} = 1$$

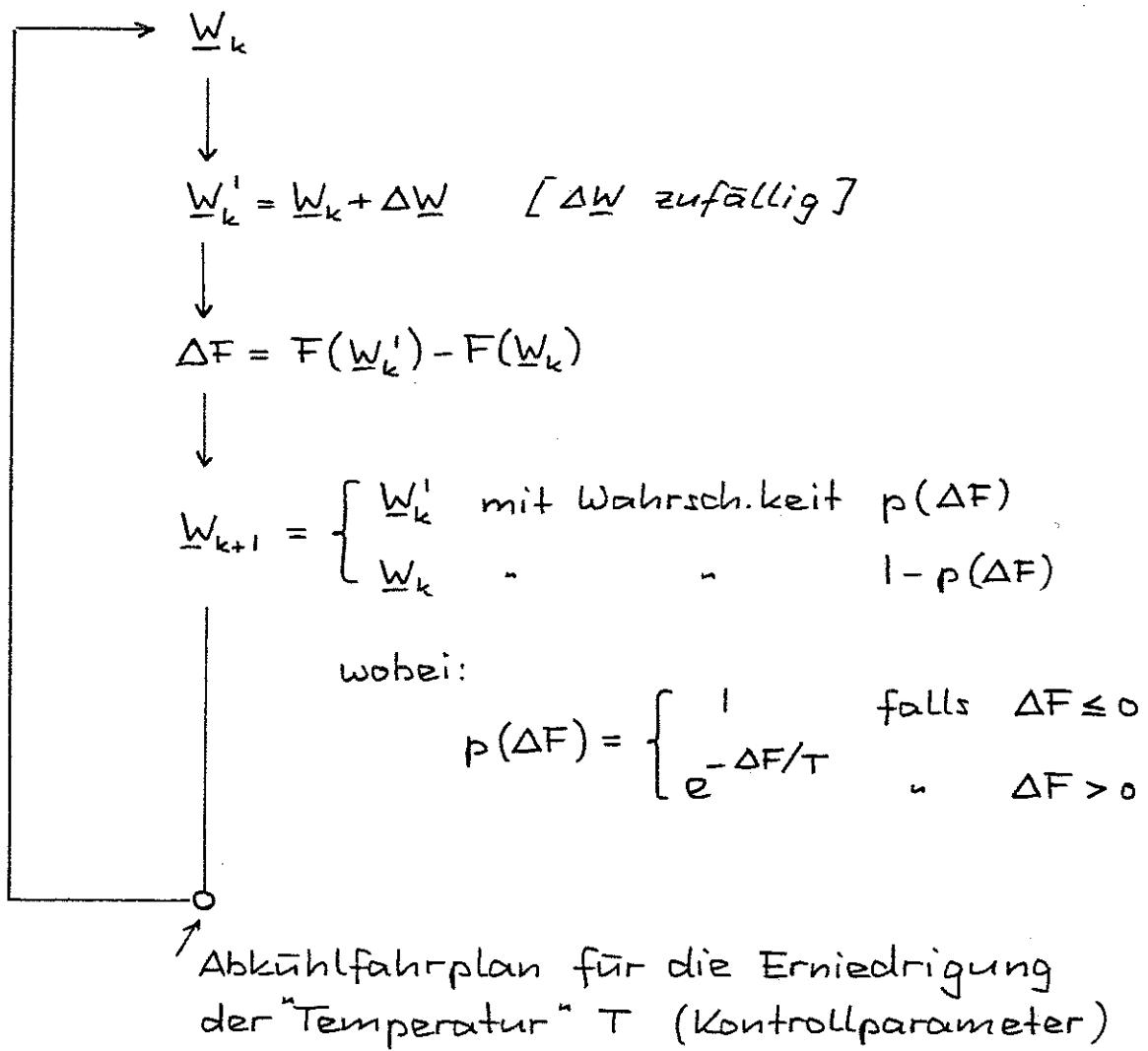
[F.J. Solis and R.J.-B. Wets, Math. of Op. Res. 6, 19 (1981)]

d.h. das "Iterative Improvement" Verfahren findet das Minimum von  $F$  im Unterraum  $S$  mit beliebiger Genauigkeit, falls genügend Iterationsschritte gemacht werden.

Frage: Verhalten des Algorithmus für beschränkte Zahl von Iterationsschritten?

## VIII.4 "Simulated Annealing"

Der Simulated Annealing Algorithmus stammt aus der statistischen Physik [Kirkpatrick et al, Science 220, 671 (1983)] und kann wie folgt beschrieben werden:



### Bemerkungen:

- Für jedes  $T > 0$  sind Gewichtsänderungen, die den Fehler vergrößern, mit einer gewissen Wahrsch.keit zugelassen. → Möglichkeit, aus schlechten lokalen Minima wieder herauszukommen.
- "Iterative Improvement" ist ein Spezialfall von "Simulated Annealing" ( $T=0$  von Anfang an).

## Warum funktioniert Simulated Annealing?

Betrachte ein physikalisches System, dessen Zustand durch  $\underline{w}$  charakterisiert ist und dessen Energiefunktion durch  $F(\underline{w})$  gegeben ist.

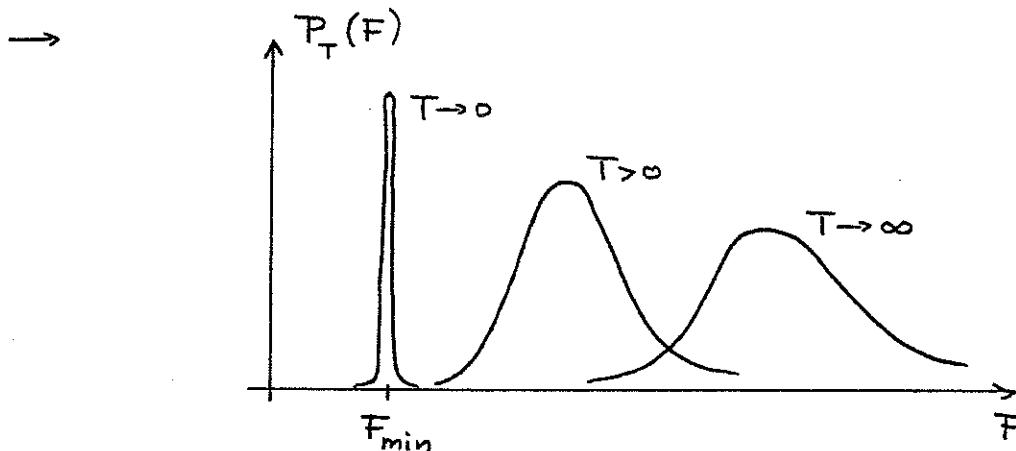
Suche Zustände  $\underline{w}$  mit möglichst tiefer "Energie"  $F(\underline{w})$ .

### Statistische Mechanik:

→ Wahrscheinlichkeit  $P_T(\underline{w})$ , das System bei einer Temperatur  $T$  im Zustand  $\underline{w}$  zu finden:

$$P_T(\underline{w}) = \frac{1}{Z_T} e^{-F(\underline{w})/T}, \quad Z_T = \sum_{\underline{w}} e^{-F(\underline{w})/T}$$

(Boltzmann-Verteilung)      (Zustandssumme)



Mittlere "Energie":  $\langle F \rangle(T) = - \frac{d \ln Z_T}{d(1/T)}$

Spezifische Wärme:  $C(T) = \frac{d \langle F \rangle}{dT} = \frac{1}{T^2} [\langle F^2 \rangle - \langle F \rangle^2]$

Entropie:  $S(T) = \frac{d}{dT} (T \ln Z_T)$

- 1) Das Annealing - Verfahren simuliert die Gleichgewichtsverteilung  $P_T(\underline{W})$ , falls es nur genügend lange bei einer konstanten Temperatur  $T$  iteriert wird. (Unter sehr allgemeinen Bedingungen an die zulässigen  $\Delta \underline{W}$ .)
- 2) Durch langsame Abkühlen verschiebt sich diese Verteilung dann zu immer tieferen  $F$ -Werten.

Optimierung des Abkühlfahrplans?

- Falls  $T_k = \frac{A}{\log(1+k)}$ ,  $k=1, 2, \dots$ ,

so konvergiert der Simulated Annealing Algorithmus gegen das globale Minimum von  $F$ :

$$\lim_{k \rightarrow \infty} F(\underline{W}_k) = F_{\min}.$$

[Mitra et al., Adv. Appl. Prob. 18, 747 (1986)]

- Optimierung des Abkühlfahrplans für vorgegebenes  $k_{\max} < \infty$  (endliche CPU-Zeit!) ??  
→ Anpassung an die "thermodynamischen Eigenschaften" des Problems!

[v. Laarhoven and Aarts:

Simulated Annealing - Theory and Applications  
Reidel Publishing (1987)]

[Bennasconi: Optimization Problems and Statistical Mechanics  
in "Chaos and Complexity", World Scientific (1988)]

## VIII.5 Beispiele für stochastische Optimierung

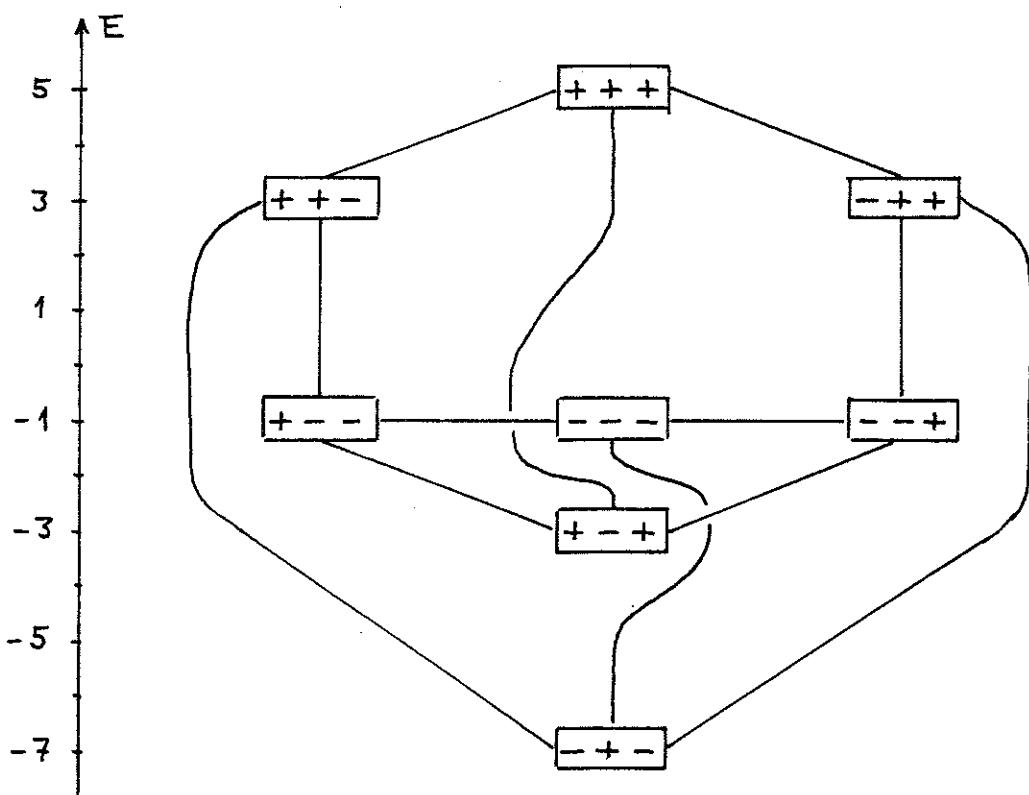
### BEISPIEL 1:

$$E = W_1 + W_2 + W_3 + 2W_1W_2 - W_1W_3 + 2W_2W_3 - W_1W_2W_3 = \min$$

wobei  $W_1, W_2$  und  $W_3$  nur die Werte +1 oder -1 annehmen dürfen

Erlaubte zufällige Änderungen:

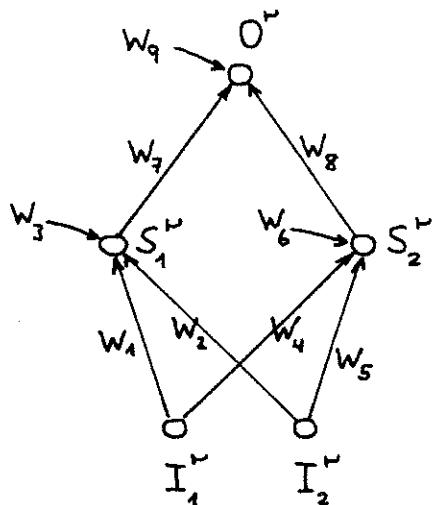
Vorzeichenwechsel eines einzelnen  $W_i$



- Das "Iterative Improvement"-Verfahren kann im lokalen Minimum  $[+ - +]$  steckenbleiben, da nur Übergänge zu Zuständen mit tieferem  $E$  zugelassen sind.
- Mit "Simulated Annealing" kann man aus diesem Minimum wieder herauskommen und in das globale Minimum  $- + -$  gelangen!

BEISPIEL 2 :

XOR-Problem, 2 Hidden Units



$$O^r = f(W_7 S_1^r + W_8 S_2^r + W_9)$$

$$S_1^r = f(W_1 I_1^r + W_2 I_2^r + W_3)$$

$$S_2^r = f(W_4 I_1^r + W_5 I_2^r + W_6)$$

$$f(x) = \frac{1 - e^{-x}}{1 + e^{-x}}$$

r	$I_1^r$	$I_2^r$	$D^r$ (XOR)
1	1	1	-1
2	1	-1	1
3	-1	1	1
4	-1	-1	-1

$$E = \frac{1}{2} \sum_{r=1}^4 (D^r - O^r)^2 = \min$$

Stochastische Optimierung:

→ Zufällige Änderungen:

$$W_i \longrightarrow W_i + \lambda \cdot \text{RND}, \quad \lambda = 0.5$$

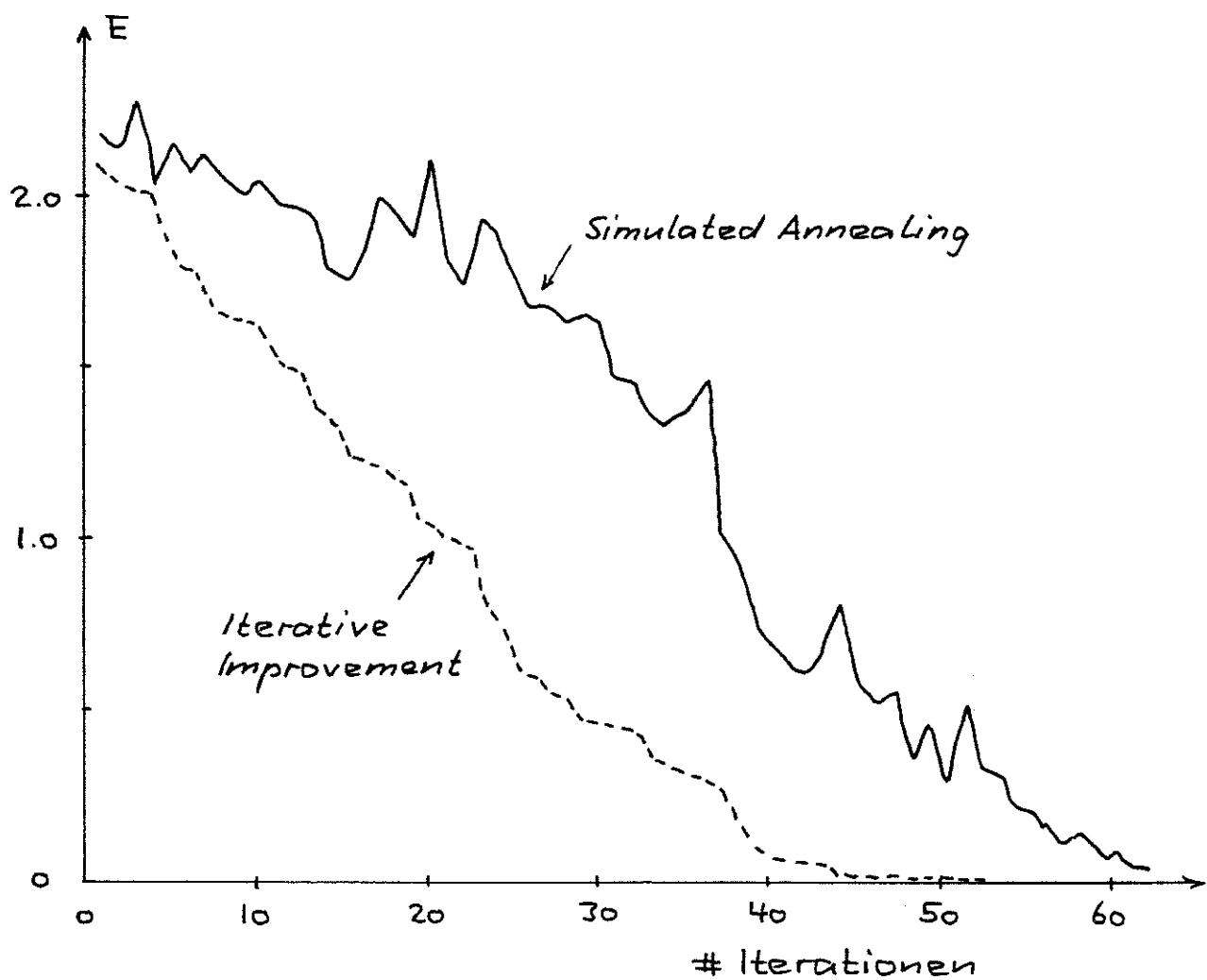
RND = random im  
Intervall (-1, +1)

[1 Iteration = alle Gewichte 1 mal zufällig geändert]

Resultate:

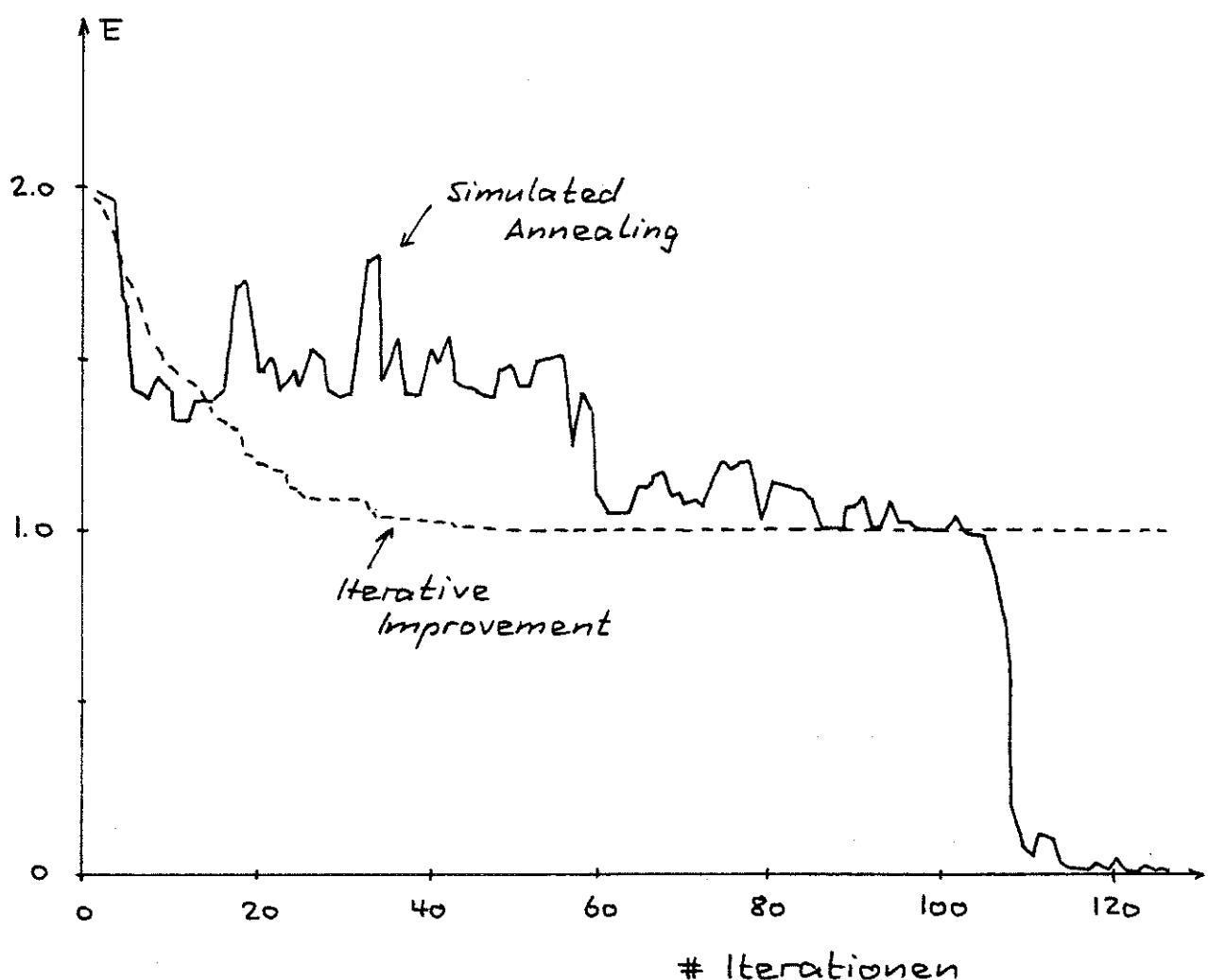
a) Wenn die Anfangsgewichte im Attraktionsgebiet des globalen Minimums liegen, konvergiert das "Iterative Improvement"-Verfahren schneller als der "Simulated Annealing"-Algorithmus.

Beim "Iterative Improvement" nimmt E immer ab, während "Simulated Annealing" auch Übergänge zu Zuständen mit höherem E zulässt.



[XOR-Problem , 2 Hidden Units]

- b) Oft bleibt das "Iterative Improvement"-Verfahren aber in einem schlechten lokalen Minimum stecken.  
 "Simulated Annealing" hingegen kann aus einem solchen Minimum wieder herauskommen. Nach genügend vielen Iterationen (theoretisch allerdings nur nach  $\infty$  vielen) findet "Simulated Annealing" daher immer das globale Minimum.



[XOR-Problem, 2 Hidden Units]

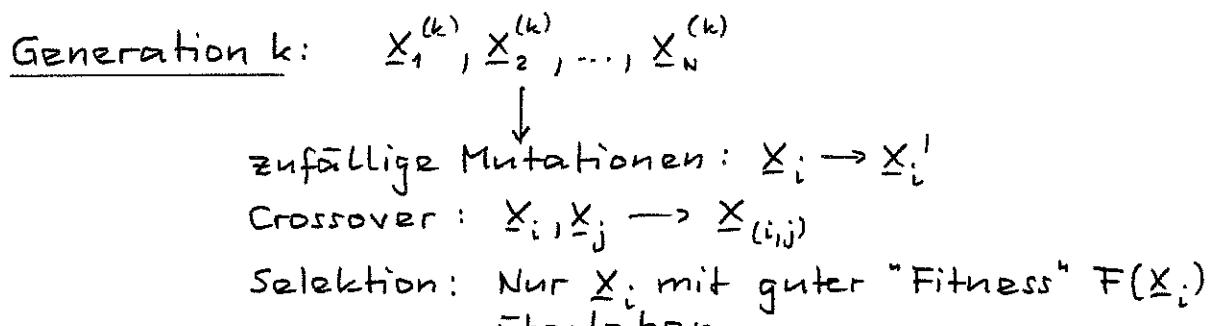
## VIII.6 Neuroevolution

Verwendung von genetischen Algorithmen (Evolutionsstrategien) zum Training von neuronalen Netzwerken.

Solche Lernverfahren operieren auf einer Population von verschiedenen neuronalen Netzwerken.

"Genotyp" eines neuronalen Netzwerks beschrieben durch Codierung der Gewichte oder/und der Netzwerk-Struktur:

- (i)  $\underline{X} = \underline{W}$  (Netzwerk-Struktur fest vorgegeben)
- (ii)  $\underline{X} = \underline{S}$  (Optimierung der Netzwerk-Struktur; Training der Gewichte z.B. mit Error-Backprop.)
- (iii)  $\underline{X} = \{\underline{S}, \underline{W}\}$  (Optimierung von Struktur und Gewichten)



Generation k+1:  $\underline{x}_1^{(k+1)}, \underline{x}_2^{(k+1)}, \dots, \underline{x}_N^{(k+1)}$

→ Population entwickelt sich zu Individuen mit optimaler "Fitness"  $F(\underline{x}_i)$

### Beispiele:

- $F(\underline{x}_i)$  = Summe der quadr. Output-Fehler für Lernbeispiele
- $F(\underline{x}_i)$  = Aufprall-Geschwindigkeit (Mondlandung)
- $F(\underline{x}_i)$  = Zeit bis zum "Umfallen" (invertiertes Pendel)
- $F(\underline{x}_i)$  = Gewonnene Punkte bei Spielen

### Literatur-Hinweise:

- G. Weiss, Towards the Synthesis of Neural and Evolutionary Learning  
in "Progress in Neural Networks" 5 (1997), 145-176
- M. Mandisch, Neurocomputing 42 (1-4), 2002, 87-117

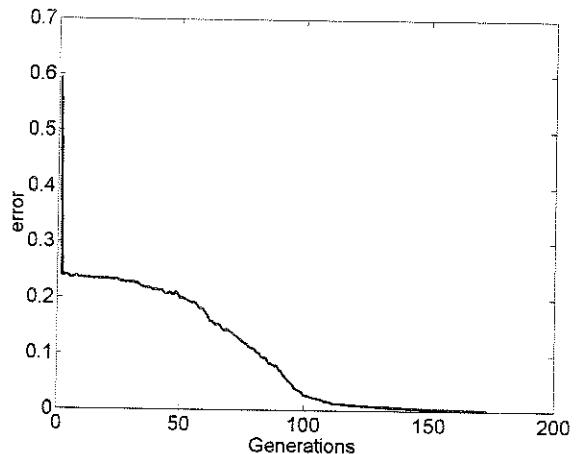
## VIII.7 Beispiele für Neuroevolution

### BEISPIEL 1: Symmetry-Detection

(Beat Gfeller, Semesterarbeit 2005)

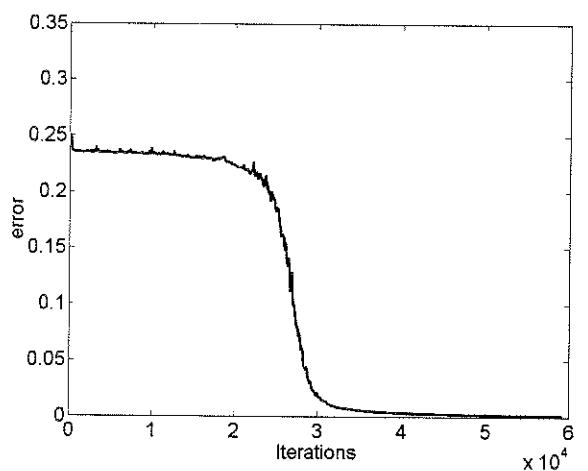
- 8-6-2-1 Multilayer-Perceptron  
→ 71 Gewichte (inkl. Schwellen)
- Output =  $\begin{cases} +1 & \text{falls binärer Input symmetrisch} \\ -1 & \text{sonst} \end{cases}$
- 15 Individuen / Generation
- 100 Netzwerk-Evaluationen / Generation

Evolutions-  
algorithmus:

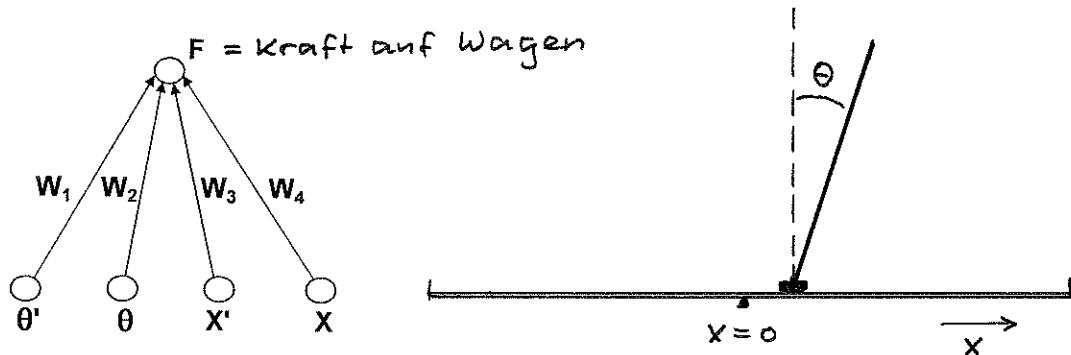


→ 17'000 Netzwerk-Evaluationen  
bis Fehler < 0.001

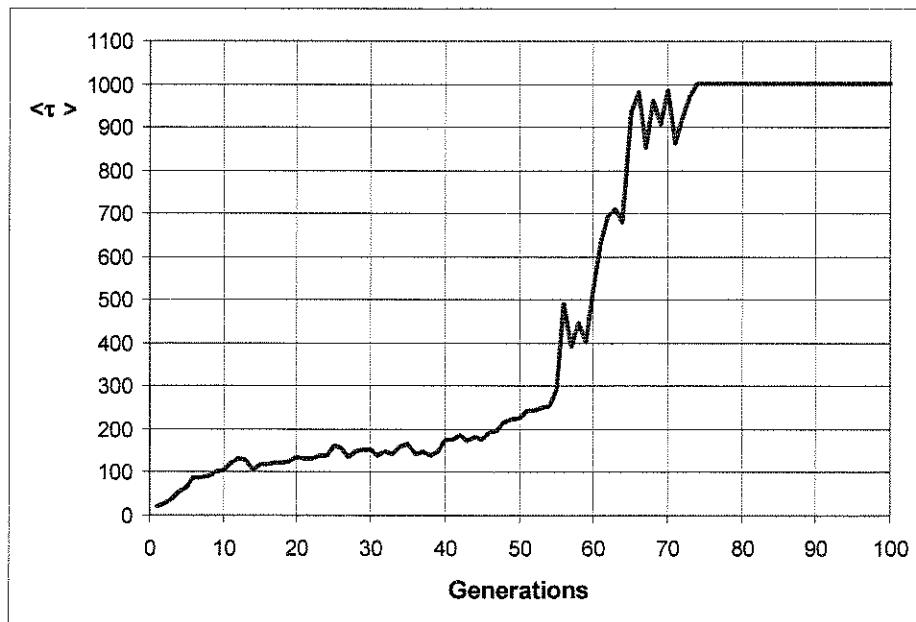
Error-  
Backpropagation:



→ 60'000 Iterationen bis Fehler < 0.001

BEISPIEL 2 : Pole-Balancing

- Individuen = Perceptrons (4 Gewichte)
- 8 Individuen / Generation
- Fitness-Funktion = "Time to Failure"



$\langle \tau \rangle$  = "Average Time to Failure"  
für 4 beste Individuen / Generation  
( $\tau = 1000$  :  $\rightarrow$  Ende der Episode)

→ ca. 70 Generationen  
bis  $\tau > 1000$  sec