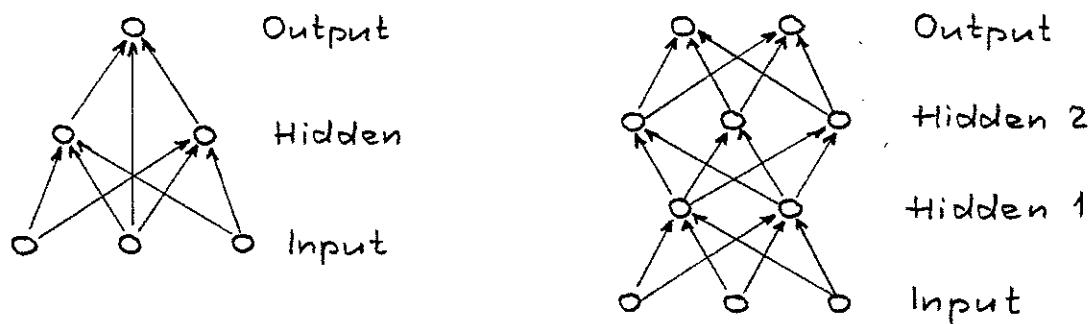


VI. MULTILAYER PERCEPTRONS

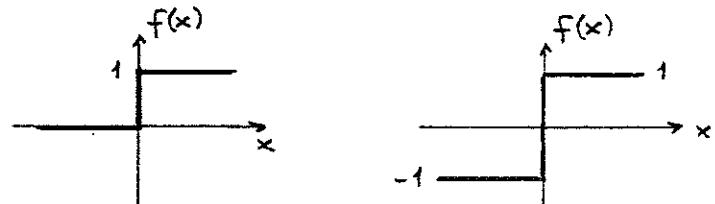


- Feedforward-Netzwerke mit einer oder mehreren Schichten von Hidden Units
- Beliebig komplizierte Klassifizierungsgebiete realisierbar
- Lernverfahren?

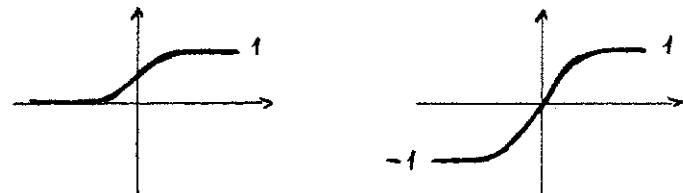
Aktivierungsfunktionen:

$$S_i \xrightarrow{W_i} \dots \xrightarrow{W_n} S = f(\sum_i W_i S_i)$$

- Stufenfunktionen:



- Sigmoidal Funktionen:



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

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

→ $f'(x) = f(x)[1-f(x)]$ $f'(x) = \frac{1}{2} [1 - f(x)^2]$

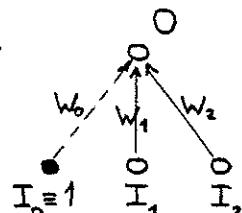
VI.1 Rolle der Hidden Units

Die "hidden units" (inneren Neuronen) konstruieren eine interne Darstellung der Inputmuster, welche es erlaubt, die gewünschte Input/Output - Abbildung zu realisieren.

Beispiel: XOR - Problem : $I_1 \quad I_2 \quad \text{XOR}$

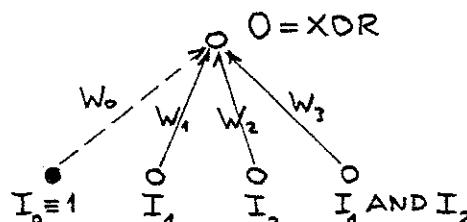
			AND
+	+	-	+
+	-	+	-
-	+	+	-
-	-	-	-

- Perceptron



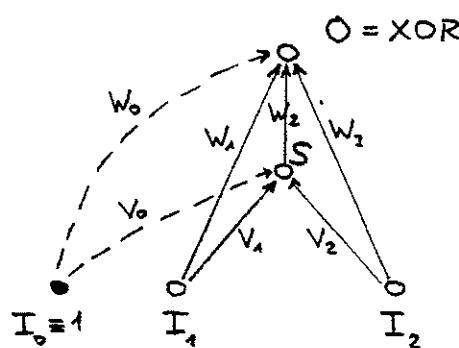
kann XOR - Abbildung nicht implementieren !

- Aber erweiterte Input - Darstellung macht Perceptron - Lösung möglich :



z.B. $w_0 = -1$
 $w_1 = +1$
 $w_2 = +1$
 $w_3 = -2$

- Realisierung mit 1 hidden Unit :



Das "Hidden" Neuron S konstruiert eine geeignete interne Darstellung (z.B. $S = I_1 \text{ AND } I_2$) automatisch während dem Lernprozess !

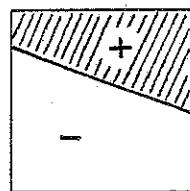
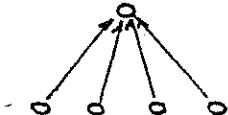
VI.2 Lernvermögen

1) Realisierbare Klassifizierungsgebiete:

- Inputs reell ($-\infty < I_i < +\infty$)
- 1 Output mit $f(x) = \text{sign}(x)$ [d.h. $O = \pm 1$]

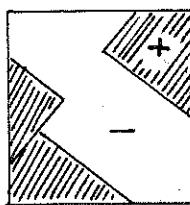
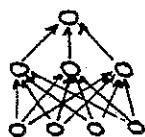
Inputraum:

Perceptron:



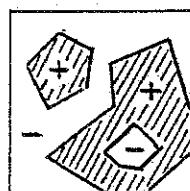
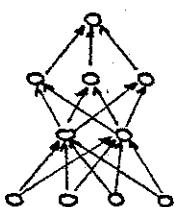
Durch Hyperebene
begrenzter Halbraum

1 hidden layer:



Nichtüberlappende
konvexe Polygone
(offen oder abgeschlossen)

2 hidden layers:



Beliebige Gebiete
(Komplexität nur
durch Zahl der Hidden
Units begrenzt)



2 Schichten von Hidden Units
genügen, um jede Klassifizierung
realisieren zu können!

2) Boole'sche Funktionen:

- Inputs I_1, \dots, I_N , $I_j = \pm 1$

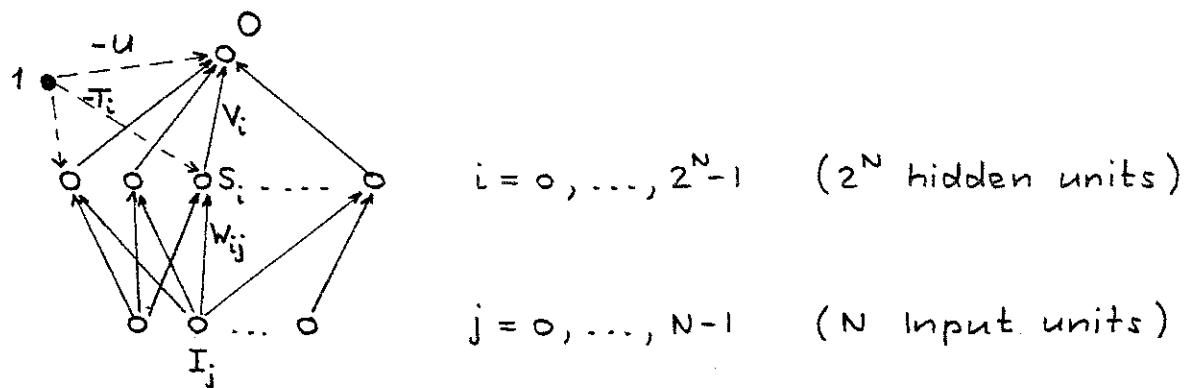
- Output $O = \pm 1$

$$\rightarrow \text{Boole'sche Funktion } O = f_B(I_1, \dots, I_N)$$



Ein Netzwerk mit 1 Schicht hidden Units kann jede beliebige Boole'sche Funktion realisieren!

Beweis konstruktiv:



$$S_i = \text{sign}(\sum_j w_{ij} I_j - T_i), \quad O = \text{sign}(\sum_i v_i S_i - u)$$

wähle: $w_{ij} = \begin{cases} +b & \text{falls "j-tes bit der Binärdarst. von } i\text{"} = 1 \\ -b & \text{sonst} \end{cases}$

$$T_i = (N-1)b$$

\rightarrow Für i -tes Inputmuster ($i=0, \dots, 2^N - 1$) ist $S_i = +1$ und $S_k = -1$, $k \neq i$.

wähle weiter:

$$V_i = p D_i \quad , \quad D_i = \text{vorgegebener Output für } i\text{-tes Inputmuster}$$

$$U = - \sum_i V_i$$

$$\rightarrow O = D_i \quad \text{für } i\text{-tes Inputmuster!}$$

Bemerkung: Der Beweis zeigt nur, dass es mit 2^n Hidden Units geht, aber nicht wieviele Hidden Units es mindestens braucht.

3) Stetige Funktionen:

→ Jede stetige Funktion $\mathbb{R}^n \rightarrow \mathbb{R}^m$ kann durch ein Feedforward-Netzwerk [n Input units, m (lineare) Output units] mit nur 1 Schicht Hidden Units [Γ oder \cup] beliebig genau approximiert werden!

[K. Hornik et al, Neural Networks 2, 359 (1989)]

ZUSAMMENFASSUNG:

- Jedes Klassifizierungs- oder Modellierungsproblem lässt sich durch ein neuronales Netzwerk mit höchstens 2 Schichten von Hidden Units lösen!

[Das bedeutet aber nicht, dass die Beschränkung auf 1 oder 2 hidden layers immer optimal ist, z.B. bezüglich Lernverhalten!]

VI.3 Lernen durch Fehlerminimierung

- Lernbeispiele: $\underline{I}^r / \underline{D}^r$, $r = 1, \dots, N$
- \underline{D}^r = vorgegebener Output
- Output des Netzwerks: $\underline{O}^r = \underline{o}(\underline{I}^r, \{W_{ij}\})$
- Quadratischer Outputfehler:

$$F_r = \frac{1}{2} \sum_i (D_i^r - O_i^r)^2 \quad [\text{Fehler f\"ur Beispiel } r]$$

$$F = \sum_r F_r \quad [\text{totaler Outputfehler}]$$

→ **LERNEN** = Bestimmen der W_{ij} , sodass
 $F = \sum_r F_r = \min$

GRADIENTENVERFAHREN:

→ Iterative Minimierung von F : $W_{ij} \rightarrow W_{ij} + \Delta W_{ij}$

a) Example-by-Example Learning:

$$\Delta W_{ij} = -\eta \frac{\partial F_r}{\partial W_{ij}} \quad \rightarrow \begin{array}{l} \text{Gewichts\"anderung direkt} \\ \text{nach jeder Pr\"äsentation} \\ \text{eines Lernbeispiels.} \end{array}$$

(Lernbeispiele in zuf\"älliger Reihenfolge!)

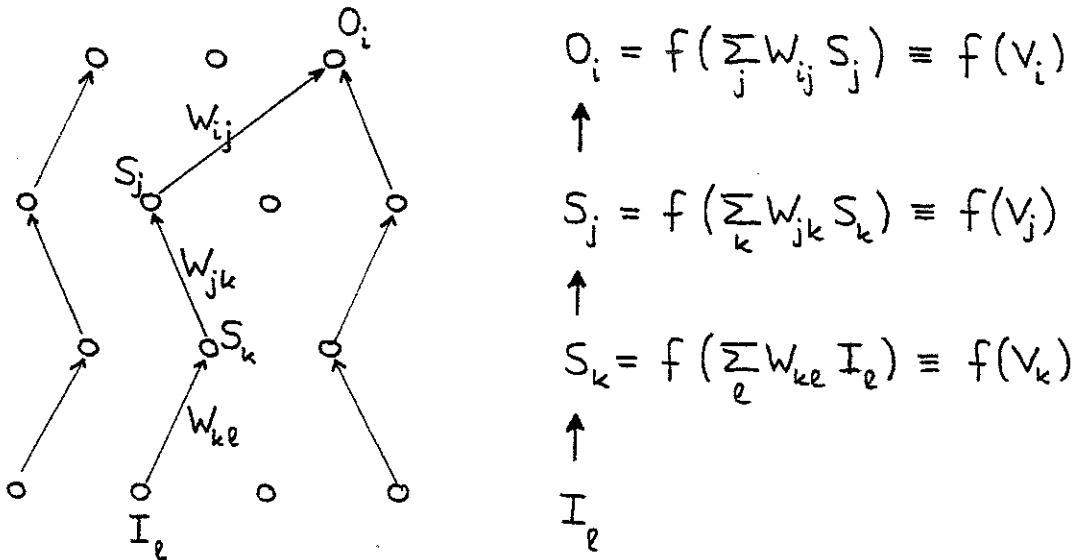
b) Batch-Learning:

$$\Delta W_{ij} = -\eta \frac{\partial F}{\partial W_{ij}}, \quad F = \sum_r F_r$$

→ Gewichts\"anderung erst nach Pr\"äsentation aller Lernbeispiele.

Bemerkung: Bei Gradientenverfahren m\"ussen Fehlerfunktion F und Aktivierungsfunktion f differenzierbar sein!

VI.4 Der "Error-Backpropagation" Algorithmus



$$F = \frac{1}{2} \sum_i (D_i - O_i)^2 = \min : \Delta W_{ij} = -\eta \frac{\partial F}{\partial W_{ij}}$$

$$\rightarrow \Delta W_{ij} = \eta (D_i - O_i) f'(V_i) S_j \equiv \delta_i \cdot S_j$$

$$\begin{aligned} \Delta W_{jk} &= \eta \sum_i (D_i - O_i) f'(V_i) W_{ij} f'(V_j) S_k \\ &= \sum_i \delta_i W_{ij} f'(V_j) S_k \equiv \delta_j \cdot S_k \end{aligned}$$

$$\Delta W_{ke} = \sum_j \delta_j W_{jk} f'(V_k) I_e \equiv \delta_k \cdot I_e$$

\rightarrow Backpropagation: $\delta_i = \eta (D_i - O_i) f'(V_i) \rightarrow \Delta W_{ij} = \delta_i S_j$

↓

$\delta_j = \sum_i \delta_i W_{ij} \cdot f'(V_j) \rightarrow \Delta W_{jk} = \delta_j S_k$

↓

$\delta_k = \sum_j \delta_j W_{jk} \cdot f'(V_k) \rightarrow \Delta W_{ke} = \delta_k I_e$

- Der Backpropagation - Algorithmus ist eine effiziente Implementierung der Gradientenmethode für Feedforward - Netzwerke.

Probleme beim Backpropagation - Lernverfahren :

Fehlerfunktion $F(\underline{w})$ ist oft sehr kompliziert und kann viele lokale Minima und lange, flache Täler enthalten:

- schlechte Lösungen (lokale Minima)
- Lernprozess sehr langsam (flache Täler, η klein)
- Oszillationen, keine Konvergenz (η zu gross)

VI.5 Varianten, "Tricks"

- 1) Überrelaxation ("momentum term"):

$$\Delta W_{ij}(k+1) = -\eta \frac{\partial F}{\partial W_{ij}} + \alpha \Delta W_{ij}(k)$$

[$\Delta W_{ij}(k)$ = Gewichtsänderung nach k-ter Lernbeispielpräsentation]

- 2) Akkumulation der Gewichtsänderungen über gewisse Anzahl n von Lernbeispielen.

- 3) Gewichtszerfall:

$$\Delta W_{ij} = -\eta \frac{\partial F}{\partial W_{ij}} - \beta W_{ij}$$

- 4) Von Zeit zu Zeit zufällige Gewichtsänderungen vornehmen.

- 5) Alternative Fehlerfunktionen benutzen:

statt $F^* = \frac{1}{2} \sum_i (D_i^* - O_i^*)^2$

(Summe der quadrierten Outputfehler)

z.B.:

$$F^* = \sum_i \left[D_i^* \log \frac{D_i^*}{O_i^*} + (1-D_i^*) \log \frac{1-D_i^*}{1-O_i^*} \right] \quad \text{falls } 0 < D, O < 1$$

$$F^* = \sum_i \left[\frac{1}{2} (1+D_i^*) \log \frac{1+D_i^*}{1+O_i^*} + \frac{1}{2} (1-D_i^*) \log \frac{1-D_i^*}{1-O_i^*} \right] \quad \text{falls } -1 < D, O < +1$$

("Informationsdifferenz"!)

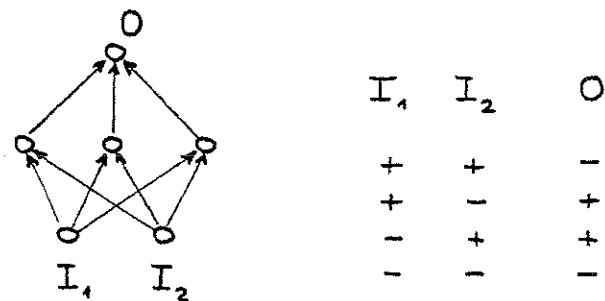
- 6) Überhaupt kein Gradientenverfahren zur Minimierung von F verwenden, sondern z.B. ein stochastisches Optimierungsverfahren:

- Simulated Annealing
- Evolutionsalgorithmen
- etc

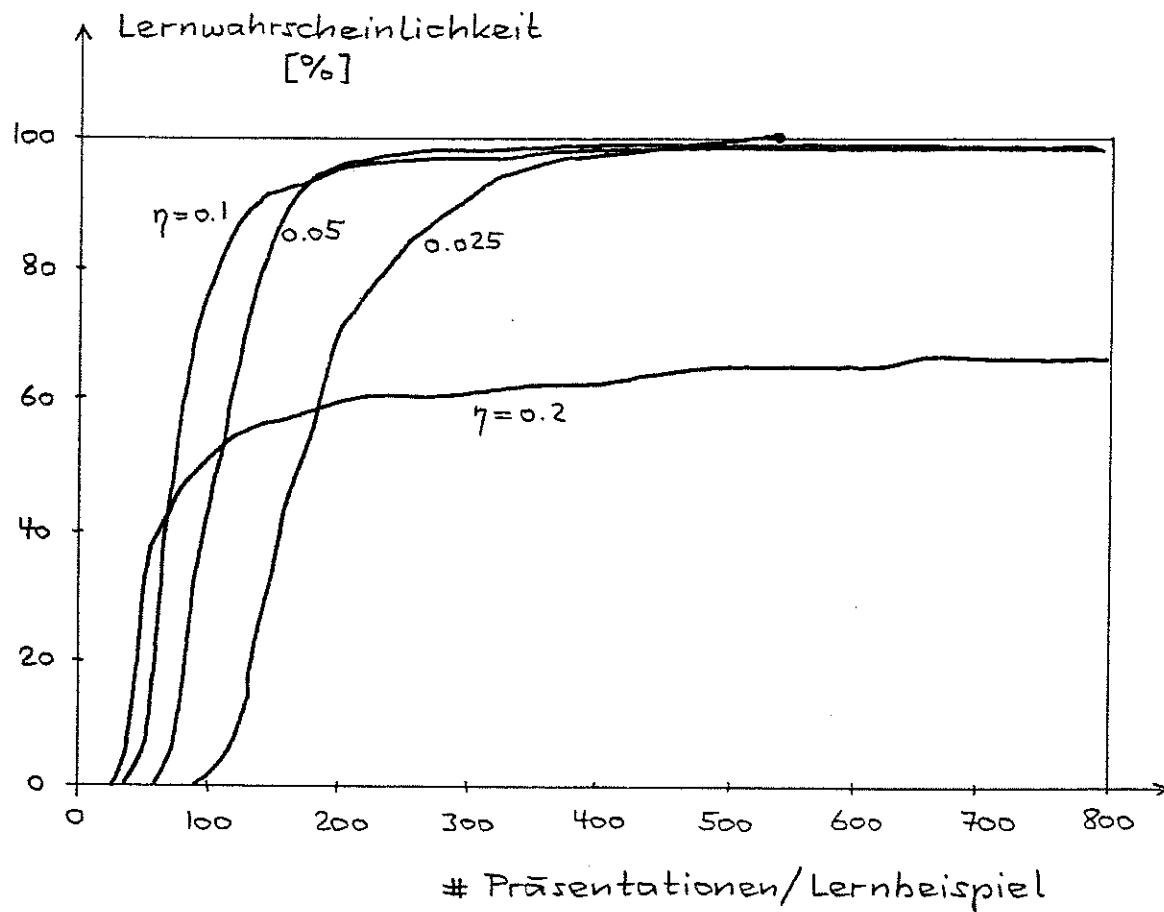
(siehe Kapitel VIII)

VI.6 Beispiele

A) XOR-Problem:

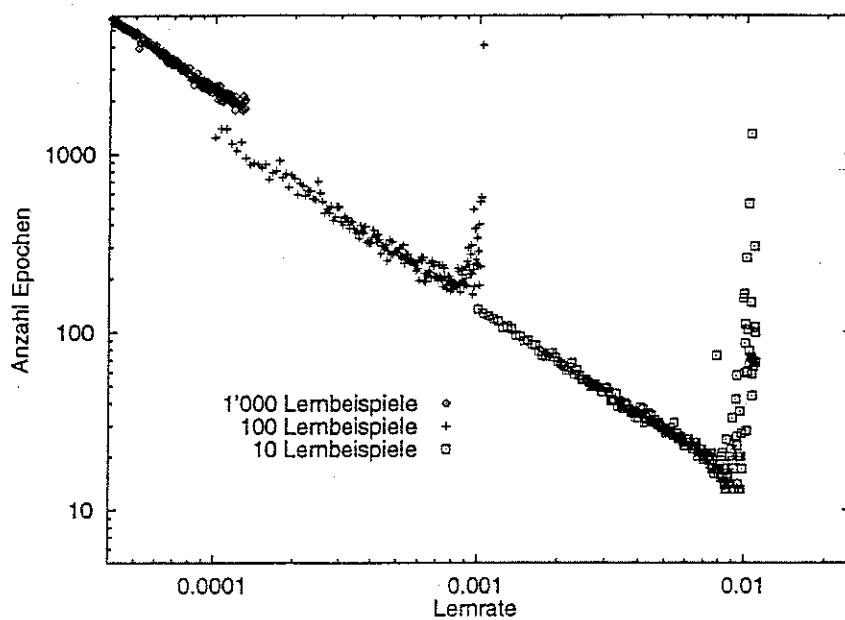


- 3 hidden units: \rightarrow 13 Gewichte (inkl. Schwellen)
- Backpropagation $[f(x) = \frac{1-e^{-x}}{1+e^{-x}}, 0 > 0.8 \text{ bzw. } < -0.8]$
- $\eta = 0.025, 0.05, 0.1, 0.2$; $\alpha = 0.9$
- Gewichtsänderung nach jeder Präsentation eines Lernbeispiels

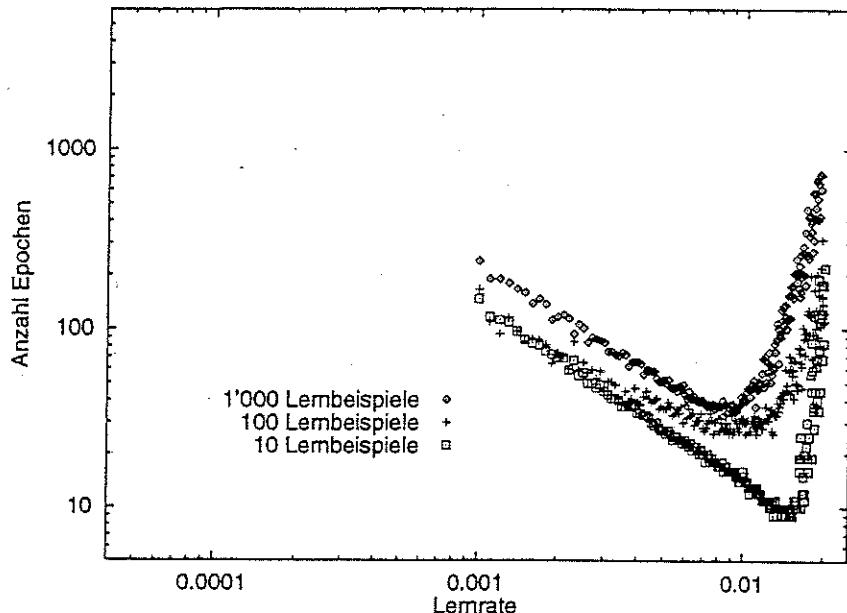


B) Erkennung von handgeschriebenen Ziffern (Urs Müller, 1995)

- Gegeben: Handgeschriebene Ziffern (0 bis 9), dargestellt als 16×16 Pixelbilder [Pixel = -1 (weiss) oder +1 (schwarz)]
- Multilayer-Perceptron mit 1 Schicht von Hidden Units:
256 Inputneuronen
100 Hidden Units
10 Outputneuronen (je 1 für jede Ziffer)
Aktivierungsfunktion:
 $f(x) = 1.1 \cdot \tanh(1.5x)$
- Es wird trainiert, bis alle Lernbeispiele richtig erkannt werden. (Error-Backpropagation)
- Figuren B1 und B2 zeigen Anzahl notwendige Epochen (Epoche = Präsentation aller Lernbeispiele) als Funktion der Lernrate η für verschiedene grosse Lernsets.



Figur B1 : Batch-Learning



Figur B2 : Example-by-Example Learning

BEMERKUNGEN :

- Konvergenzgeschwindigkeit nimmt mit zunehmender Lernrate zuerst zu (Anzahl notwendiger Epochen sinkt) und dann wieder ab.
Bei noch grösseren Lernraten als in den Figuren eingezeichnet: \rightarrow keine Konvergenz mehr.
- Für kleine Lernsets (10 Lernbeispiele) ist die Konvergenzgeschwindigkeit für beide Verfahren (Batch und Example-by-Example) etwa gleich. Auch die optimalen Lernraten sind ähnlich.
- Für grosse Lernsets ist Example-by-Example Learning wesentlich schneller als Batch Learning.
- Batch Learning kann aber auch Vorteile haben:
z.B. werden bei der Approximation von kontinuierlichen Input/Output - Zusammenhängen mit Batch Learning grössere Genauigkeiten erzielt als mit Example-by-Example Learning!

GRADIENTEN - LERNVERFAHREN
 (ERGÄNZUNG)

NOTATION: $\underline{W} = (W_1, W_2, \dots, W_n)$, $O^\mu = O(I^\mu, \underline{W})$

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

$$\left[\text{z.B. } F(\underline{W}) = \frac{1}{2} \sum_{\mu=1}^N (D^\mu - O^\mu)^2 \right]$$

GRADIENTENVERFAHREN:

Iterationsschritte : $k = 1, 2, \dots$

$$\rightarrow \underline{W}^{k+1} = \underline{W}^k + \Delta \underline{W}^k \quad , \quad \Delta \underline{W}^k = \Delta \underline{W}(\underline{g}^k)$$

$$\underline{g} = \left(\frac{\partial F}{\partial W_1}, \dots, \frac{\partial F}{\partial W_n} \right) \quad , \quad \underline{g}^k = \underline{g}(\underline{W}^k)$$

EINFACHE GRADIENTENVERFAHREN:

[z.B. Backpropagation (Batch-Mode!)]

\underline{W}^0 zufällig

$$\underline{W}^{k+1} = \underline{W}^k - \eta \underline{g}^k$$

$\eta = \text{const.}$ (genügend klein!)

\rightarrow Probleme bei langen flachen Tälern, etc.

METHODE DER KONJUGIERTEN GRADIENTEN

[siehe z.B.:

R. Battiti, Neural Computation 4, 141-166 (1992)]

\underline{W}^0 zufällig, $\underline{P}^0 = \underline{g}^0$

$\rightarrow \underline{W}^{k+1} = \underline{W}^k - \alpha_k \underline{P}^k$, α_k so, dass F entlang \underline{P}^k minimal wird

$$\underline{P}^{k+1} = -\underline{g}^{k+1} + \beta_k \underline{P}^k, \quad \beta_k = \frac{(\underline{g}^{k+1} - \underline{g}^k) \cdot \underline{g}^{k+1}}{(\underline{g}^{k+1} - \underline{g}^k) \cdot \underline{P}^k}$$

- Falls F = quadr. Funktion der W_{ij} :

\rightarrow Auf der Geraden durch \underline{W}^{k+1} in Richtung \underline{P}^{k+1} ist der Gradient von F in Richtung \underline{P}^k immer gleich Null [k-te Minimierung wird nicht zerstört!].

\rightarrow Verfahren konvergiert in n Schritten!
(n = Anzahl Gewichte)

- Neuronale Netzwerke:

$\rightarrow F$ ist keine quadr. Funktion der W_{ij}

\rightarrow Numerische Probleme

\rightarrow Varianten von KG-Verfahren

Für viele Lernprobleme sind KG-Methoden wesentlich schneller als Backpropagation!