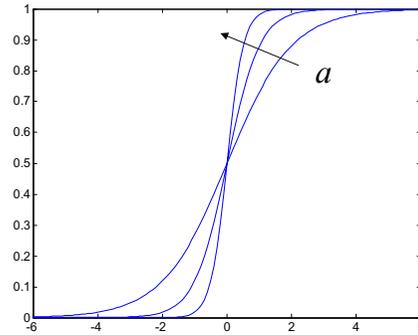


Lernstrategien für Neuronale Netze - Der Backpropagation-Algorithmus

Im Unterschied zu den bisher behandelten NNs mit Schwellwertfunktionen sind die NNs mit *stetig differenzierbaren nichtlinearen Aktivierungsfunktionen* $f(x)$. Am weitesten verbreitet ist die *Sigmoid-Funktion* $\psi(x)$:

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

Sie approximiert mit wachsendem a die Sprungfunktion $\sigma(x)$.



Die Funktion $\psi(x)$ hat einen engen Bezug zur tanh-Funktion. Es gilt für $a=1$:

$$\tanh(x) = \frac{2}{1 + e^{-2x}} - 1 = 2\psi(2x) - 1$$

Fügt man die zwei Parameter a und b hinzu, so erhält man eine etwas allgemeinere Form:

$$\psi_g(x) = \frac{1}{1 + e^{-a(x-b)}} - 1$$

wobei mit a die Steilheit und mit b die Position auf der x -Achse beeinflusst werden kann.

Die Nichtlinearität der Aktivierungsfunktion ist entscheidend für die Existenz des Multi-Lagen-Perceptrons; ohne diese, würde das mehrschichtige in ein triviales lineares einschichtiges Netzwerk zusammenschrumpfen.

Die Differenzierbarkeit von $f(x)$ erlaubt die Anwendung von notwendigen Bedingungen ($\partial(\cdot)/\partial w_{ij}$) zur Optimierung der Gewichtskoeffizienten w_{ij} .

Die erste Ableitung der Sigmoid-Funktion kann auf sich selbst zurückgeführt werden:

$$\frac{d\psi}{dx} = \frac{-1}{(1 + e^{-x})^2} (-1)e^{-x} = \frac{1}{1 + e^{-x}} \left(1 - \frac{1}{1 + e^{-x}}\right) = \psi(x)(1 - \psi(x))$$

Eine Schicht des Neuronalen Netzes

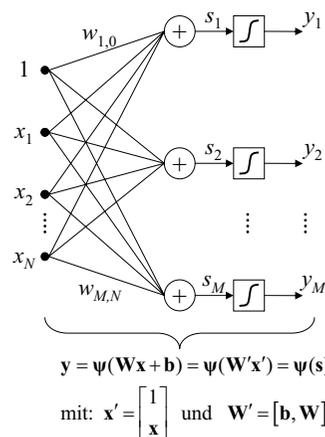
Eine einzige Schicht des NN wird charakterisiert durch die $M \times N$ Koeffizienten der Gewichtsmatrix \mathbf{W} , den Offset-Vektor \mathbf{b} und eine vektorielle Sigmoid-Funktion ψ .

Nach einer Erweiterung des Eingangsvektors um eine konstante 1

$$\mathbf{x}' = \begin{bmatrix} 1 \\ \mathbf{x} \end{bmatrix}$$

ergibt sich eine Schicht in der nebenstehenden Form.

Ihre Funktion lässt sich beschreiben durch eine Matrixmultiplikation, gefolgt von einer nichtlinearen Aktivierungsfunktion, welche in identischer Form auf alle Elemente angewandt wird.

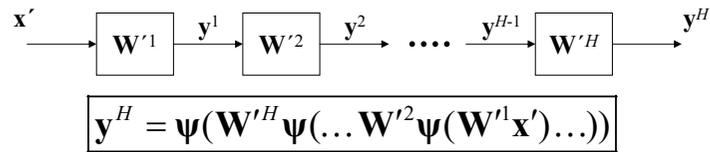


Gestalt der erweiterten Gewichtsmatrix

$$\mathbf{W}' = \begin{bmatrix} w_{1,0} = b_1 & w_{1,1} & w_{1,2} & \cdots & w_{1,N} \\ w_{2,0} = b_2 & w_{2,1} & w_{2,2} & \cdots & w_{2,N} \\ w_{3,0} = b_3 & w_{3,1} & w_{3,2} & \cdots & w_{3,N} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ w_{M,0} = b_M & w_{M,1} & w_{M,2} & \cdots & w_{M,N} \end{bmatrix} = [\mathbf{b}, \mathbf{W}]$$

Das Neuronale Netz mit H Schichten

Das mehrlagige NN mit H Schichten hat in jeder Schicht eine eigene Gewichtsmatrix \mathbf{W}^i , jedoch identische Sigmoid-Funktionen:



Das Lernen basiert auf der Anpassung der Gewichtsmatrizen, mit dem Ziel der Minimierung eines Fehlerquadrat-Kriteriums zwischen dem Sollwert \mathbf{y} und der Approximation $\hat{\mathbf{y}}$ durch das Netz (überwachtes Lernen). Der Erwartungswert ist zu bilden über das zur Verfügung stehende Trainings-Ensemble von $\{\hat{\mathbf{y}}_j, \mathbf{x}_j\}$:

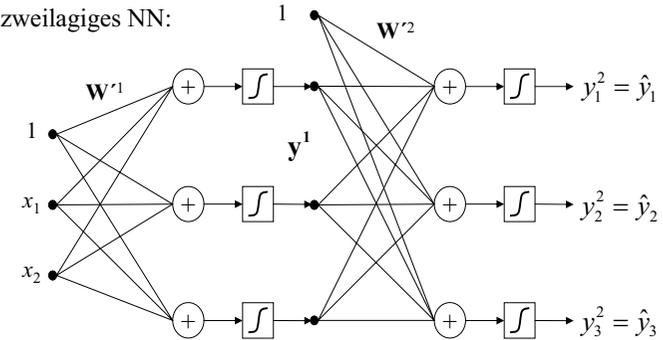
$$J = \min_{\mathbf{w}^i} E \left\{ \|\hat{\mathbf{y}} - \mathbf{y}\|^2 \right\} = \min_{\mathbf{w}^i} E \left\{ \|\mathbf{y}^H - \mathbf{y}\|^2 \right\}$$

Anmerkung: Das einschichtige NN und der lineare Polynomklassifikator sind identisch, wenn die Aktivierungsfunktion $\psi \equiv 1$ gesetzt wird.

Beziehungen zu dem Konzept der Funktionsapproximation durch eine Linearkombination von Basisfunktionen

Bei der Approximation z.Bsp. mit Polynomfunktionen sind die Basisfunktionen von vornherein festgelegt. Bei dem NN werden die Basisfunktionen iterativ durch Parametermatrizen aufgebaut. Zudem ist die Linearkombination nicht der letzte Schritt der Approximation; vielmehr geht die Wirkung durch eine nichtlineare Abbildung hervor und der Vorgang wiederholt sich!

Beispiel für ein zweilagiges NN:

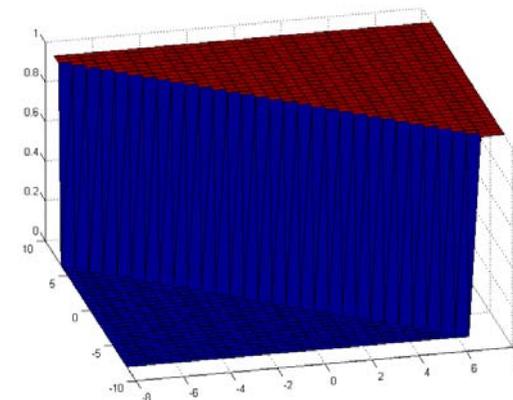


Beziehungen zu dem Konzept der Funktionsapproximation durch eine Linearkombination von Basisfunktionen

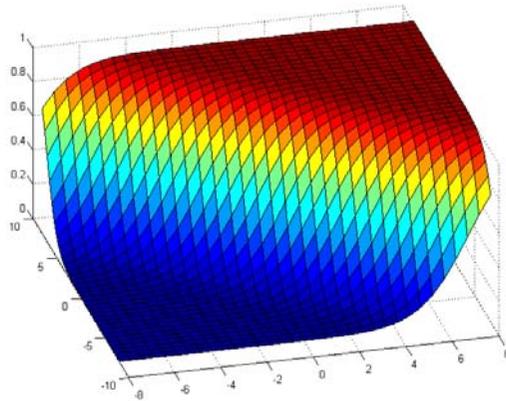
Die erste Schicht erzeugt die Basisfunktionen in Form des verdeckten Vektors \mathbf{y}_1 und die zweite Schicht bildet eine Linearkombination dieser Basisfunktionen.

Deshalb beeinflusst die Koeffizientenmatrix \mathbf{W}^1 der ersten Schicht das *Aussehen* der Basisfunktionen, während die Gewichtsmatrix \mathbf{W}^2 der zweiten Schicht die *Koeffizienten der Linearkombinationen* enthält. Diese wird zusätzlich durch die Aktivierungsfunktion gewichtet.

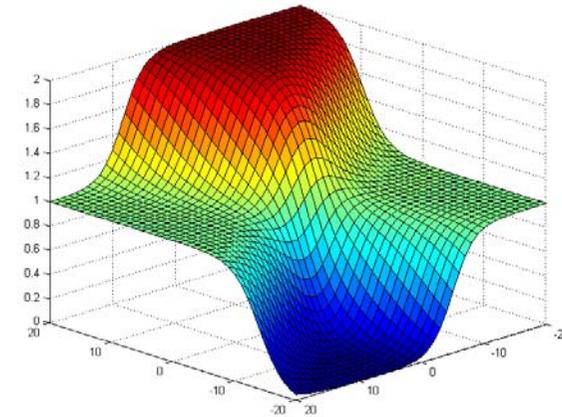
Reaktion eines Neurons mit zwei Eingängen und einer Schwellwertfunktion σ



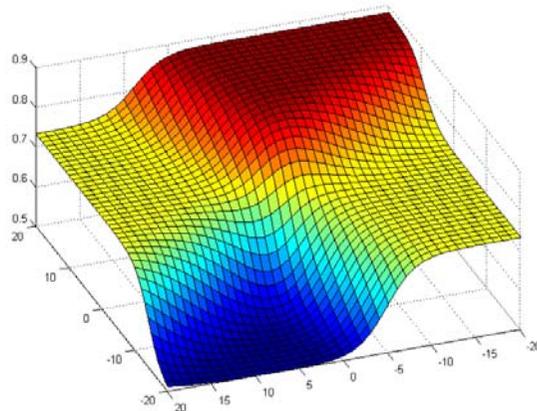
Reaktion eines Neurons mit zwei Eingängen und einer Sigmoidfunktion ψ



Überlagerung zweier Neuronen mit je zwei Eingängen und einer Sigmoidfunktion



Reaktion eines Neurons in der zweiten Schicht auf zwei Neuronen der ersten Schicht nach der Bewertung durch die Sigmoidfunktion



Der Backpropagation-Lernalgorithmus

- Die Klassenzugehörigkeitsabbildung geschieht durch ein Multilagenperceptron, welches am Ausgang eine 1 erzeugt in den Regionen von \mathbf{x} , welche durch die Stichproben \mathbf{x}_i der entsprechenden Bedeutungsklasse bevölkert ist, und sie erzeugt eine 0 in Gebieten, welche durch andere Bedeutungsklassen belegt sind. In den Gebieten dazwischen und außerhalb findet eine Inter- bzw. eine Extrapolation statt.
- Das Netzwerk wird trainiert auf der Grundlage der Optimierung des quadratischen Gütekriteriums:

$$J = E \left\{ \left\| \hat{\mathbf{y}}(\mathbf{W}^i, \mathbf{x}') - \mathbf{y} \right\|^2 \right\}$$

- Dieser Ausdruck ist nichtlinear sowohl in den Elementen des Eingangsvektors \mathbf{x}' , als auch in den Gewichtskoeffizienten $\{w'_{ij}\}$.

- Nach Einsatz der Funktionsabbildung des Multilagenperceptrons erhält man:

$$J = E \left\| \underbrace{\Psi(\mathbf{W}^{rH} \dots \Psi(\mathbf{W}^{r2} \Psi(\mathbf{W}^{r1} \mathbf{x}')) \dots)}_{\hat{\mathbf{y}}} - \mathbf{y} \right\|^2$$

- Gesucht ist das *globale* Minimum. Es ist jedoch nicht klar, ob ein solches existiert oder wie viele *lokale* Minima vorhanden sind. Der Backpropagation Algorithmus löst das Problem iterativ mit einem Gradientenalgorithmus. Iteriert wird gemäß:

$$\mathbf{W}' \leftarrow \mathbf{W}' - \alpha \nabla \mathbf{J} \quad \text{mit: } \mathbf{W}' = \{\mathbf{W}^{r1}, \mathbf{W}^{r2}, \dots, \mathbf{W}^{rH}\}$$

$$\text{und d. Gradienten: } \nabla \mathbf{J} = \frac{\partial J}{\partial \mathbf{W}'} = \left\{ \frac{\partial J}{\partial w_{nm}^{rh}} \right\}$$

- Die Iteration wird beendet, wenn der Gradient bei einem (lokalen oder auch globalen) Minimum verschwindet.
- Eine geeignete Wahl von α ist schwierig. Kleine Werte erhöhen die Anzahl der Iterationen. Größere Werte vermindern zwar die Wahrscheinlichkeit in ein lokales Minimum zu laufen, aber man riskiert, daß das Verfahren divergiert und das Minimum nicht gefunden wird (oder Oszillationen auftauchen).

- Die Iteration hat leider nur eine *lineare* Konvergenzordnung (Gradientenalgorithmus).

Berechnung des Gradienten:

Zur Bestimmung der zusammengesetzten Funktion

$$\hat{\mathbf{y}}(\mathbf{x}') = \Psi(\mathbf{W}^{rH} \dots \Psi(\mathbf{W}^{r2} \Psi(\mathbf{W}^{r1} \mathbf{x}')) \dots)$$

- wird die Kettenregel benötigt.
- Der Fehler, verursacht durch eine Stichprobe ergibt sich zu:

$$J_j = \frac{1}{2} \|\hat{\mathbf{y}}(\mathbf{x}_j) - \mathbf{y}_j\|^2 = \frac{1}{2} \sum_{k=1}^{N^H} (\hat{y}_k(j) - y_k(j))^2$$

- Der Erwartungswert $E\{\dots\}$ des Gradienten muss er durch den Mittelwert über alle Stichprobengradienten approximiert werden:

$$\nabla \mathbf{J} = \frac{1}{n} \sum_{j=1}^n \nabla \mathbf{J}_j$$

- Man unterscheidet zwischen
 - individuellem Lernen aufbauend auf die letzte Stichprobe

$$[\mathbf{x}_j, \mathbf{y}_j] \Rightarrow \nabla \mathbf{J}_j$$

- und dem kumulativen Lernen (batch learning), aufbauend auf

$$\nabla \mathbf{J} = \frac{1}{n} \sum_{j=1}^n \nabla \mathbf{J}_j$$

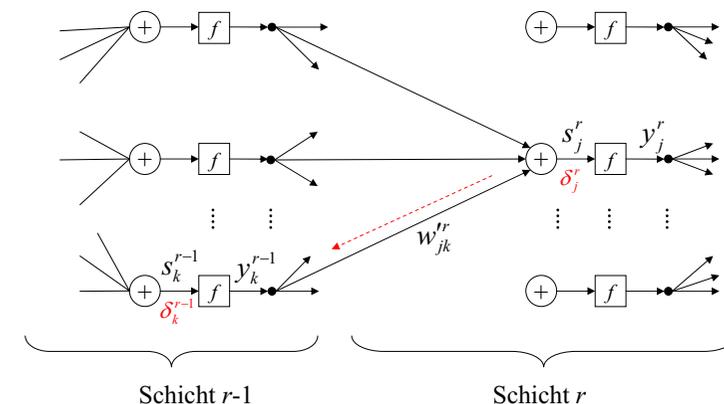
- Partielle Ableitungen für die eine Schicht:
Für die r -te Schicht gilt:

$$y_m^r = \Psi(s_m^r) = \Psi \left(\sum_{n=0}^{N^r} w_{nm}^{rr} x_n^{rr} \right)$$

x_n^{rr} n-ter Eingang von Schicht r

y_m^r m-ter Ausgang von Schicht r

Definition der Variablen zweier Schichten für den Backpropagation-Algorithmus



- Wir berechnen zunächst die Wirkung der letzten verdeckten Schicht auf die Ausgangsschicht $r=H$.
Unter Verwendung der partiellen Ableitungen der Gütefunktion J (eigentlich J_j , aber j wird der Einfachheit halber weggelassen) und unter Anwendung der Kettenregel erhält man:

$$\frac{\partial J}{\partial w_{nm}^{rH}} = \frac{\partial J}{\partial s_n^H} \cdot \frac{\partial s_n^H}{\partial w_{nm}^{rH}} = -\delta_n^H \underbrace{\frac{\partial s_n^H}{\partial w_{nm}^{rH}}}_{=y_m^{r-1}}$$

- unter Einführung der *Empfindlichkeit* von Zelle n

$$\delta_n = -\frac{\partial J}{\partial s_n}$$

- ergibt sich:
$$\delta_n^H = -\frac{\partial J}{\partial s_n^H} = -\frac{\partial J}{\partial \hat{y}_n^H} \cdot \frac{\partial \hat{y}_n^H}{\partial s_n^H} = \underbrace{(y_n - \hat{y}_n)}_{\partial \left(\frac{1}{2} \sum_{k=1}^{N^H} (\hat{y}_k - y_k)^2 \right)} f'(s_n^H)$$
- und damit schließlich für das Update der Gewichte:

$$w_{neu} = w_{alt} + \Delta w \quad \text{mit} \quad \Delta w = -\alpha \frac{\partial J}{\partial w}$$

$$\Delta w_{nm}^{rH} = -\alpha \delta_n^H y_m^{r-1} = -\alpha (y_n - \hat{y}_n) f'(s_n^H) y_m^{r-1}$$

- Für alle anderen verdeckten Schichten $r < H$ sind die Überlegungen etwas komplexer. Wegen der Abhängigkeit der Schichten untereinander, beeinflussen die Werte von s_j^{r-1} alle Elemente s_k^r der nachfolgenden Schicht. Unter Anwendung der Kettenregel erhält man:

$$\underbrace{\frac{\partial J}{\partial s_j^{r-1}}}_{\delta_j^{r-1}} = \sum_k \underbrace{\frac{\partial J}{\partial s_k^r}}_{\delta_k^r} \frac{\partial s_k^r}{\partial s_j^{r-1}}$$

- Mit

$$\frac{\partial s_k^r}{\partial s_j^{r-1}} = \frac{\partial \left(\sum_m w_{km}^{rr} \overbrace{y_m^{r-1}}^{f(s_m^{r-1})} \right)}{\partial s_j^{r-1}} = w_{kj}^{rr} f'(s_j^{r-1})$$

- ergibt sich:

$$\delta_j^{r-1} = f'(s_j^{r-1}) \left[\sum_k w_{kj}^{rr} \delta_k^r \right]$$

- D.h. die Fehler „backpropagieren“ von der Ausgangs- zu niederen Schichten!

- Man durchläuft alle Lernstichproben, ohne irgendwelche Veränderungen an den Gewichtskoeffizienten, und man erhält den Gradienten $\nabla \mathbf{J}$ durch Mittelung über die $\nabla \mathbf{J}_j$. Beide Größen können zum Aufdaten der Parametermatrix \mathbf{W}' des Perceptrons verwendet werden:

$$\mathbf{W}'_{k+1} = \mathbf{W}'_k - \alpha \nabla \mathbf{J}_j \quad \text{individuelles Lernen}$$

$$\mathbf{W}'_{k+1} = \mathbf{W}'_k - \alpha \nabla \mathbf{J} \quad \text{kumulatives Lernen}$$

- Die Gradienten $\nabla \mathbf{J}$ und $\nabla \mathbf{J}_j$ unterscheiden sich. Der zweite ist der Mittelwert des ersten ($\nabla \mathbf{J} = E\{\nabla \mathbf{J}_j\}$), oder: der erste stellt einen zufälligen Wert des zweiten dar.

Die gesamte individuelle Lernprozedur besteht aus den folgenden Verarbeitungsschritten:

- Wähle vorläufige Werte für die Koeffizientenmatrizen $\mathbf{W}^1, \mathbf{W}^2, \dots, \mathbf{W}^H$ aller Schichten
- Eine neue Beobachtung $[\mathbf{x}'_j, \hat{\mathbf{y}}_j]$ sei gegeben
- Berechne die Diskriminierungsfunktion $\hat{\mathbf{y}}_j$ aus der gegebenen Beobachtung \mathbf{x}'_j und den momentanen Gewichtsmatrizen (Vorwärtsrechnung)
- Berechne den Fehler zwischen Schätzung $\hat{\mathbf{y}}$ und Zielvektor \mathbf{y} :
$$\Delta \mathbf{y} = \hat{\mathbf{y}} - \mathbf{y}$$
- Berechne den Gradienten $\nabla \mathbf{J}$ bzgl. aller Perceptron-Gewichte (error backpropagation). Berechne dazu zuerst δ_j^H der Ausgangsschicht gemäß:

$$\delta_j^H = -\frac{\partial J}{\partial s_n^H} = -\frac{\partial J}{\partial \hat{y}_n^H} \cdot \frac{\partial \hat{y}_n^H}{\partial s_n^H} = (y_j - \hat{y}_j) f'(s_j^H)$$

- und berechne daraus rückwärts alle Werte der niederen Schichten gemäß:

$$\delta_j^{r-1} = f'(s_j^{r-1}) \left[\sum_k w_{kj}^{rr} \delta_k^r \right] \quad \text{für } r = H, H-1, \dots, 2$$

- unter Beachtung der ersten Ableitung der Sigmoidfunktion $f(s) = \psi(s)$:

$$f'(s) = \frac{d\psi(s)}{ds} = \psi(s)(1 - \psi(s)) = y(1 - y)$$

- Korrigiere (parallel) alle Perceptron-Gewichte gemäß:

$$w_{nm}^{rr} \leftarrow w_{nm}^{rr} - \alpha \frac{\partial J}{\partial w_{nm}^{rr}} = w_{nm}^{rr} - \alpha \delta_n^r y_m^{r-1} \quad \text{für } \begin{matrix} r = 1, 2, \dots, H \\ n = 1, 2, \dots, N^H \\ m = 1, 2, \dots, M^H \end{matrix}$$

- Beim individuellen Lernen werden die Gewichte $\{w_{nm}^{r,h}\}$ mit $\nabla \mathbf{J}$ für jede Stichprobe korrigiert, wohingegen beim kumulativen Lernen die gesamte Lernstichprobe durchgearbeitet werden muss, um den gemittelten Gradienten $\nabla \mathbf{J}$ aus der Sequenz der $\{\nabla \mathbf{J}\}$ zu ermitteln, bevor die Gewichte korrigiert werden können gemäß:

$$w_{nm}^{rr} \leftarrow w_{nm}^{rr} - \sum_i \alpha \delta_n^r(i) y_m^{r-1}(i) \quad \text{für } \begin{matrix} r = 1, 2, \dots, H \\ n = 1, 2, \dots, N^H \\ m = 1, 2, \dots, M^H \end{matrix}$$

Backpropagation-Algorithmus in Matrixschreibweise

- Wähle vorläufige Werte für die Koeffizientenmatrizen $\mathbf{W}^1, \mathbf{W}^2, \dots, \mathbf{W}^H$ aller Schichten
- Eine neue Beobachtung $[\mathbf{x}'_j, \hat{\mathbf{y}}_j]$ sei gegeben
- Berechne die Diskriminierungsfunktion $\hat{\mathbf{y}}_j$ aus der gegebenen Beobachtung \mathbf{x}'_j und den momentanen Gewichtsmatrizen (Vorwärtsrechnung). Speichere alle Werte von \mathbf{y}^r und \mathbf{s}^r in allen Zwischenschichten $r = 1, 2, \dots, H$.
- Berechne den Fehler zwischen Schätzung $\hat{\mathbf{y}}$ und Zielvektor \mathbf{y} am Ausgang:
 $\Delta \mathbf{y} = \hat{\mathbf{y}} - \mathbf{y}$
- Berechne den Gradienten $\nabla \mathbf{J}$ bzgl. aller Perceptron-Gewichte (error backpropagation). Berechne dazu zuerst δ^H der Ausgangsschicht gemäß:

$$\delta^H = -\frac{\partial J}{\partial \mathbf{s}^H} = -\frac{\partial J}{\partial \hat{\mathbf{y}}^H} \cdot \frac{\partial \hat{\mathbf{y}}^H}{\partial \mathbf{s}^H} = (\mathbf{y} - \hat{\mathbf{y}}) f'(\mathbf{s}^H)$$

- und berechne daraus rückwärts alle Werte der niederen Schichten gemäß:

$$\delta^{r-1} = f'(\mathbf{s}^{r-1}) \circ (\mathbf{W}^{rT} \delta^r) \quad \text{für } r = H, H-1, \dots, 2$$

- unter Beachtung der ersten Ableitung der Sigmoidfunktion $f(s) = \psi(s)$:

$$f'(s) = \frac{d\psi(s)}{ds} = \psi(s)(1 - \psi(s)) = y(1 - y)$$

- Individuelles Lernen:** Korrigiere (parallel) alle Perceptron-Gewichtsmatrizen mit $\nabla \mathbf{J}$ für jede Stichprobe gemäß:

$$\mathbf{W}^{rr} \leftarrow \mathbf{W}^{rr} - \alpha \frac{\partial J}{\partial \mathbf{W}^{rr}} = \mathbf{W}^{rr} - \alpha \delta^r \mathbf{y}^{(r-1)T} \quad \text{für } r = 1, 2, \dots, H$$

- Kumulatives Lernen:** es muss die gesamte Lernstichprobe durchgearbeitet werden, um den gemittelten Gradienten $\nabla \mathbf{J}$ aus der Sequenz der $\{\nabla \mathbf{J}\}$ zu ermitteln, bevor die Gewichte korrigiert werden können gemäß:

$$\mathbf{W}^{rr} \leftarrow \mathbf{W}^{rr} - \alpha \sum_j \delta_j^r \mathbf{y}_j^{(r-1)T} \quad \text{für } r = 1, 2, \dots, H$$

Eigenschaften:

- Der Backpropagation-Algorithmus ist einfach zu implementieren, aber sehr rechenaufwendig, insbesondere wenn die Koeffizientenmatrix groß ist, was zur Folge hat, dass auch die Lernstichprobe entsprechend groß sein muss. Nachteilig ist weiterhin die Abhängigkeit des Verfahrens von den Startwerten der Gewichte, dem Korrekturfaktor α und die Reihenfolge, in der die Stichproben abgearbeitet werden.
- Wie beim rekursiven Trainieren des Polynomklassifikators bleiben lineare Abhängigkeiten in den Merkmalen unberücksichtigt.
- Positiv zu vermerken ist, dass der Gradientenalgorithmus auch für sehr große Probleme verwendet werden kann.
- Die Dimension der Koeffizientenmatrix \mathbf{W} ergibt sich aus den Dimensionen des Eingangsvektors, der verdeckten Schichten, sowie des Ausgangsvektors $(N, N^1, \dots, N^{H-1}, K)$ zu:

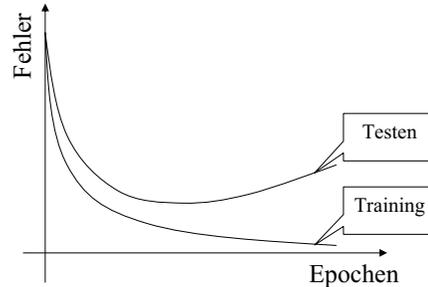
$$T = \dim(\mathbf{W}) = \sum_{h=1}^H (N^{h-1} + 1) N^h \quad \text{mit } N^0 = N \quad \text{und } N^H = K$$

$$N = \dim(\mathbf{x}) \quad \text{Merkmalsraum}$$

$$K = \dim(\hat{\mathbf{y}}) \quad \text{Anzahl der Klassen}$$

Zur Dimensionierung des Netzes

- Die Überlegungen bei dem Entwurf eines mehrschichtigen Perceptrons mit Schwellwertfunktionen (σ bzw. sign) geben eine gute Vorstellung, wieviele Hidden-Layer und wieviele Neuronen man für ein MLP verwenden sollte für das Backpropagation-Lernen (vorausgesetzt, man hat eine gewisse Vorstellung über die Verteilung der Cluster).
- Das Netz sollte so einfach wie nur möglich gewählt werden. Mit höherer Dimension steigt die Gefahr des overfitting, gepaart mit einem Verlust an Generalisierungsfähigkeit und zugleich werden viele Nebenmaxima zugelassen, wo der Algorithmus hängen bleiben kann!
- Bei einem vorgegebenen Stichprobenumfang, sollte die Lernphase bei einem bestimmten Punkt abgebrochen werden. Danach wird das Netz zu sehr an die vorhandenen Daten angepasst (overfitting) und verliert an Generalisierungsfähigkeit.



Zur Berechnungskomplexität des Backpropagation-Algorithmus

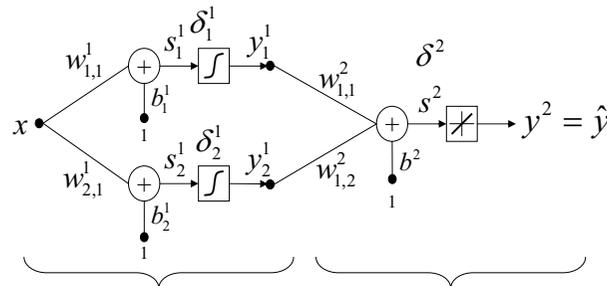
- Wenn $T = \dim(\mathbf{W})$ die Anzahl der Gewichte und Biasterme ist, so lässt sich einfach nachvollziehen, daß $O(T)$ Berechnungsschritte für die Vorwärtssimulation benötigt werden, $O(T)$ für das Backpropagieren des Fehlers und ebenso $O(T)$ Operationen für die Korrektur der Gewichte, also erhält man insgesamt eine lineare Komplexität in der Anzahl der Gewichte: $O(T)$.
- Würde man den Gradienten durch finite Differenzen experimentell bestimmen (dazu ändert man jedes Gewicht inkrementell und ermittelt die Wirkung auf das Gütemaß durch Vorwärtsrechnung) durch Berechnung eines Differenzenquotienten gemäß (d.h. man macht sich nicht die Mühe der analytischen Auswertung):

$$\frac{\partial J}{\partial w_{ji}^r} = \frac{J(w_{ji}^r + \varepsilon) - J(w_{ji}^r)}{\varepsilon} + O(\varepsilon)$$

- so ergäben sich T Vorwärtsrechnungen der Komplexität $O(T)$ und damit insgesamt eine Gesamtkomplexität von: $O(T^2)$

Backpropagation zum Trainieren eines zweischichtigen Netzwerks zur Funktionsapproximation

(Matlab-Demo: „Backpropagation Calculation“)



$$y^1 = \psi(\mathbf{W}^1 x + \mathbf{b}^1) = \psi(\mathbf{W}^1 \mathbf{x}^1) = \psi(\mathbf{s}^1) \quad \hat{y} = y^2 = (\mathbf{W}^2 \mathbf{y}^1 + b^2)$$

mit: $\mathbf{W}^1 = \begin{bmatrix} w_{1,1}^1 \\ w_{2,1}^1 \end{bmatrix}$ und $\mathbf{b}^1 = \begin{bmatrix} b_1^1 \\ b_2^1 \end{bmatrix}$ mit: $\mathbf{W}^2 = \begin{bmatrix} w_{1,2}^2 & w_{2,2}^2 \end{bmatrix}$

- Gesucht wird eine Approximation der Funktion

$$y(x) = 1 + \sin\left(\frac{\pi}{4} x\right) \quad \text{für} \quad -2 \leq x \leq 2$$

Backpropagation zum Trainieren eines zweischichtigen Netzwerks zur Funktionsapproximation

- Backpropagation:** Für die Ausgangsschicht oder zweite Schicht mit linearer Aktivierungsfunktion ergibt sich wegen $f' = 1$:

$$\delta^2 = (y - \hat{y})$$

- Backpropagation:** Für die erste Schicht mit Sigmoid-Funktion ergibt sich: :

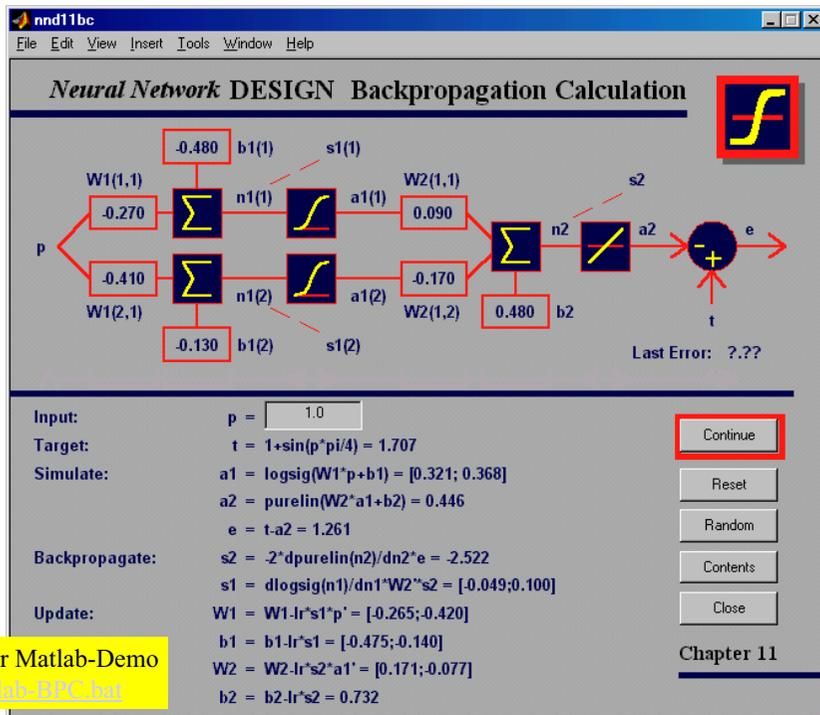
$$\delta_j^1 = f'(s_j^1) [w_{1,j}^2 \delta^2] = y_j^1 (1 - y_j^1) [w_{1,j}^2 \delta^2] \quad \text{für } j = 1, 2$$

- Korrektur der Gewichte in der Ausgangsschicht:

$$\Delta w_{1,j}^2 = \alpha \delta^2 y_j^1 \quad \text{und} \quad \Delta b^2 = \alpha \delta^2 \quad \text{für } j = 1, 2$$

- Korrektur der Gewichte in der Eingangsschicht:

$$\Delta w_{j,1}^1 = \alpha \delta_j^1 x \quad \text{und} \quad \Delta b_j^1 = \alpha \delta_j^1 \quad \text{für } j = 1, 2$$



Start der Matlab-Demo
[matlab-BPC.bat](#)

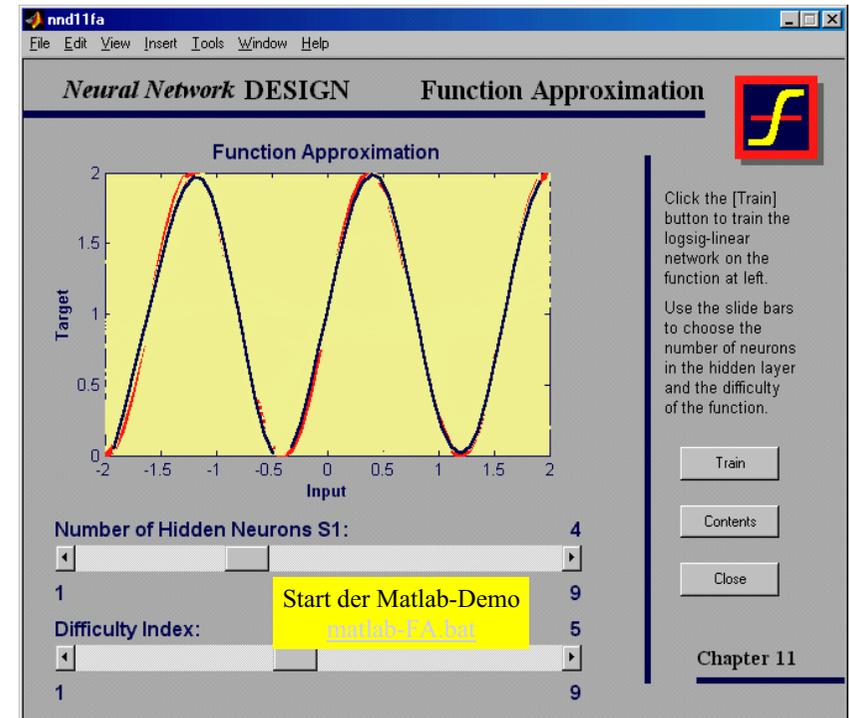
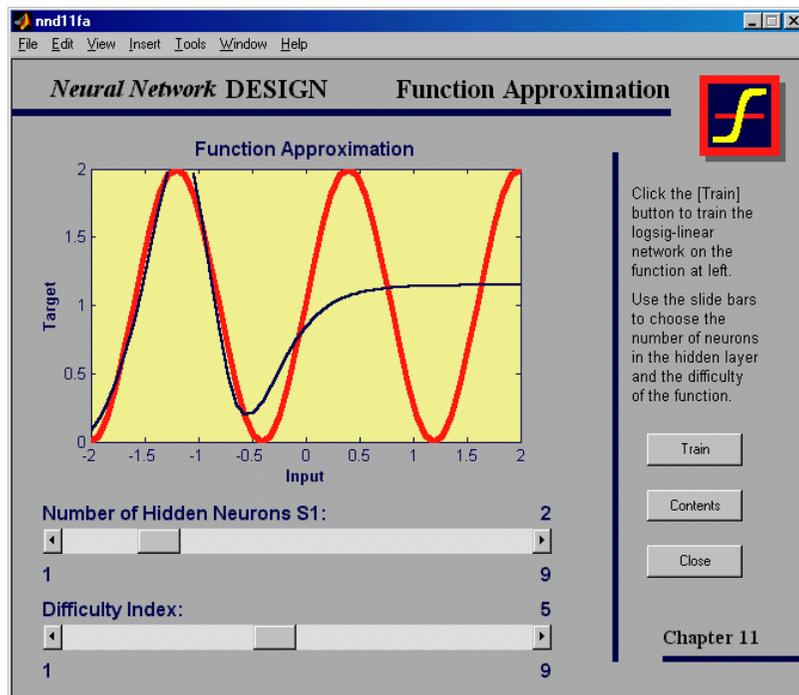
Backpropagation zum Trainieren eines zweischichtigen Netzwerks zur Funktionsapproximation

(Matlab-Demo: „Function Approximation“)

- Gesucht wird eine Approximation der Funktion

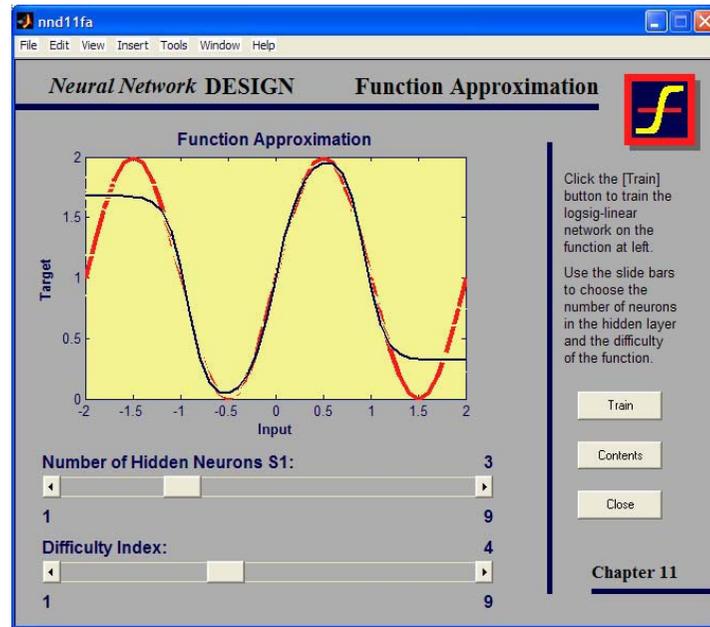
$$y(x) = 1 + \sin(i \cdot \frac{\pi}{4} x) \quad \text{für } -2 \leq x \leq 2$$

- Mit zunehmenden Wert von i (difficulty index) steigen die Anforderungen an das MLP-Netzwerk. Die auf der Sigmoid-Funktion aufbauende Approximation benötigt mehr und mehr Schichten, um mehrere Perioden der Sinus-Funktion darzustellen.
- Problem: Konvergenz zu lokalen Minima, selbst in Fällen, wo das Netzwerk groß genug gewählt wird, um die Funktion
 - Das Netzwerk wird zu klein gewählt, so dass eine Approximation nur schlecht gelingt, d.h. es wird zwar das globale Minimum erreicht, aber dieses liefert noch keine gute Approximationsgüte ($i=8$ und Netzwerk 1-3-1)
 - Das Netzwerk wird groß genug gewählt, so dass eine gute Approximation möglich ist, aber es konvergiert nur gegen ein lokales Minimum ($i=4$ und Netzwerk 1-3-1)



Start der Matlab-Demo
[matlab-FA.bat](#)

Modell groß genug, aber nur lokales Minimum



Generalisierungsfähigkeit des Netzwerkes

- Wir gehen davon aus, dass das Netzwerk für 11 Abtastpunkte trainiert wird

$$\{y_1, x_1\}, \{y_2, x_2\}, \dots, \{y_{11}, x_{11}\}$$

- Die Frage ist nun, wie gut das Netzwerk die Funktion für nicht gelernte Abtastpunkte approximiert in Abhängigkeit von der Komplexität des Netzwerkes

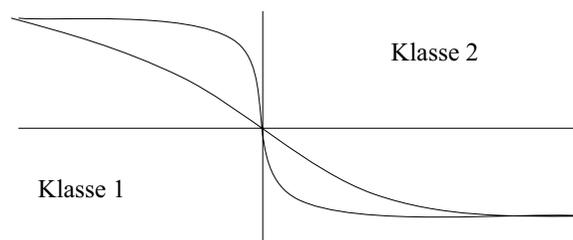
Start der Matlab-Demo (Hagan 11-21)

[matlab-GENERAL.bat](#)

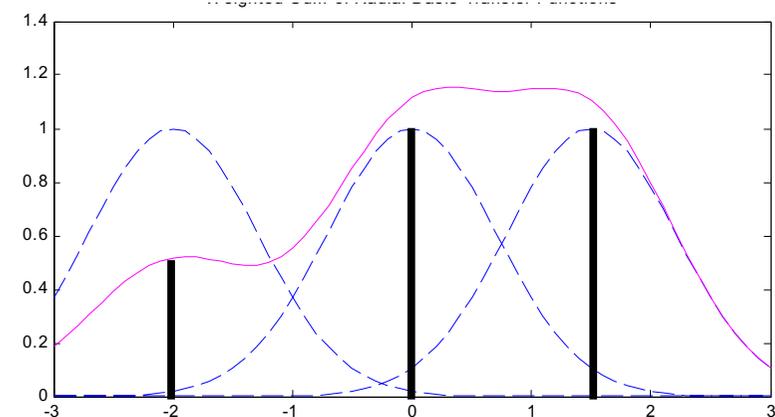
- Regel: Das Netzwerk sollte weniger Parameter haben als die zur Verfügung stehenden Ein-/Ausgangspaare

Verbesserung der Generalisierungsfähigkeit eines Netzes durch Hinzufügen additiver Störungen

- Stehen nur wenige Stichproben zur Verfügung, so kann die Generalisierungsfähigkeit eines NN verbessert werden, indem man den Stichprobenumfang durch z.Bsp. normalverteilte Störungen verbreitert. D.h. man fügt weitere Stichproben hinzu, welche mit hoher Wahrscheinlichkeit in der unmittelbaren Nachbarschaft anzutreffen sind.
- Dadurch werden die Intraklassenbereiche verbreitert und die Grenzen zwischen den Klassen schärfer ausgebildet.

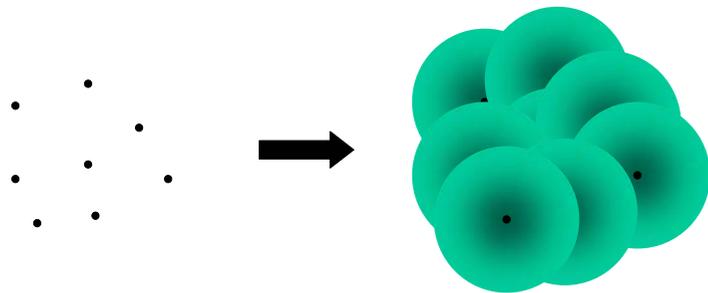


Interpolation durch Überlagerung von Normalverteilungen

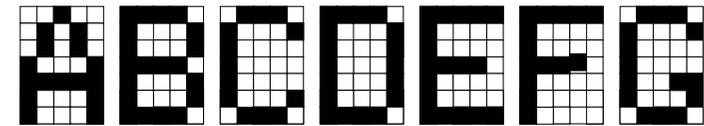


$$y(x) = \sum \alpha_i f_i(x - t_i)$$

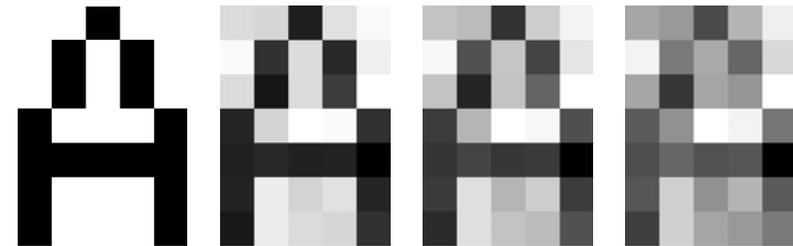
Interpolation durch Überlagerung von Normalverteilungen



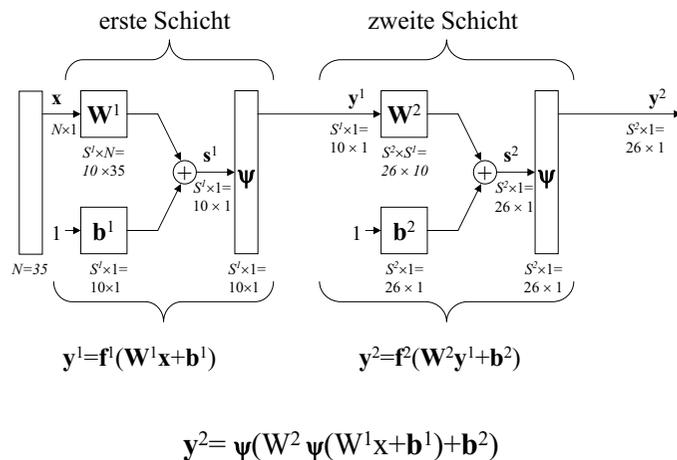
Zeichenerkennungsaufgabe für 26 Buchstaben der Größe 7x5



- mit graduell ansteigenden additiven Störungen:



Zweischichtiges Neuronales Netz für die Zeichenerkennung mit 35 Eingangswerten (Pixel), 10 Neuronen in der verdeckten Schicht und 26 Neuronen in der Ausgangsschicht



Demo mit MATLAB

- Öffnen von Matlab
 - Demo in Toolbox Neural Networks
 - „Character recognition“ (command line)
 - Es werden zwei Netzwerke trainiert
 - Netzwerk 1 ohne Störungen
 - Netzwerk 2 mit Störungen
 - Auf einem unabhängigen Testdatensatz liefert Netzwerk 2 bessere Ergebnisse als Netzwerk 1

Start der Matlab-Demo
matlab-CR.bat

Möglichkeiten zur Beschleunigung der Adaption

Es handelt sich bei der Adaption von NN um ein allgemeines Parameteroptimierungsproblem. Demgemäß können im Prinzip eine Vielzahl weiterer Optimierungsmethoden aus der numerischen Mathematik eingesetzt werden. Dies kann zu erheblichen Verbesserungen führen (Konvergenzgeschwindigkeit, Konvergenzbereich, Stabilität, Berechnungsaufwand).

I.a. lassen sich jedoch nur sehr schwer allgemeingültige Aussagen machen, da die Ergebnisse von Fall zu Fall sehr verschieden sein können und in der Regel Kompromisse in den Eigenschaften zu akzeptieren sind!

- Heuristische Verbesserungen des Gradientenalgorithmus (Schrittweitenkontrolle)
 - Gradientenalgorithmus (Steepest descent) mit Momentum-Term (update der Gewichte nicht nur abhängig vom Gradient, sondern auch vom vorherigen update)
 - Verwendung eines adaptiven Lernfaktors
- Konjugierter Gradientenalgorithmus

- Newton und Quasi-Newton-Algorithmen (Verwendung der zweiten Ableitungen der Fehlerfunktion z.B. in Form der Hesse-Matrix oder deren Schätzung)

↗ - Quickprop

✗ - Levenberg-Marquardt-Algorithmus

Taylorentwicklung der Gütefunktion in \mathbf{w} :

$$J(\mathbf{w} + \Delta\mathbf{w}) = J(\mathbf{w}) + \nabla\mathbf{J}^T \Delta\mathbf{w} + \frac{1}{2} \Delta\mathbf{w}^T \mathbf{H} \Delta\mathbf{w} + \text{Terme höherer Ordnung}$$

mit: $\nabla\mathbf{J} = \frac{\partial J}{\partial \mathbf{w}}$ Gradientenvektor

und: $\mathbf{H} = \left\{ \frac{\partial^2 J}{\partial w_i \partial w_j} \right\}$ Hesse-Matrix

- Pruning-Techniken. Man startet mit einem hinreichend großen Netzwerk und nimmt dann wieder Neuronen aus dem Netz, welche keinen oder nur geringen Einfluß auf die Gütefunktion haben und reduziert damit das Overfitting der Daten.

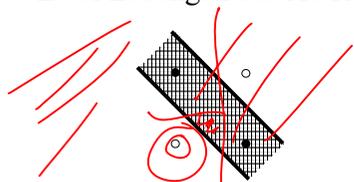
Demo mit MATLAB

(Demo von Theodoridis)

- C:\...\matlab\PR\startdemo
- Example 2 (XOR) mit (2-2-1 und 2-5-5-1)
- Example 3 mit dreischichtigem Netzwerk [5,5] (3000 Epochen, learning rate 0,7, momentum 0,5)

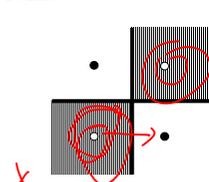
Start der Matlab-Demo
[matlab-theodoridis.bat](#)

Zwei Lösungen des XOR-Problems:



Konvexes Gebiet realisiert mit einem zweischichtigen Netz (2-2-1). Mögliche Fehlklassifikation bei Varianz:

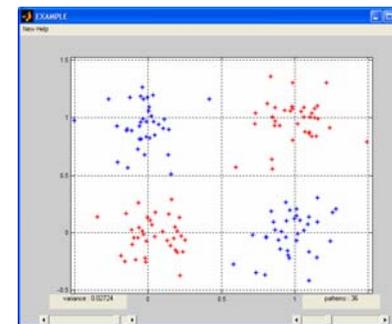
$$\|\mathbf{n}\| > \frac{1}{4}\sqrt{2} \approx 0,35$$



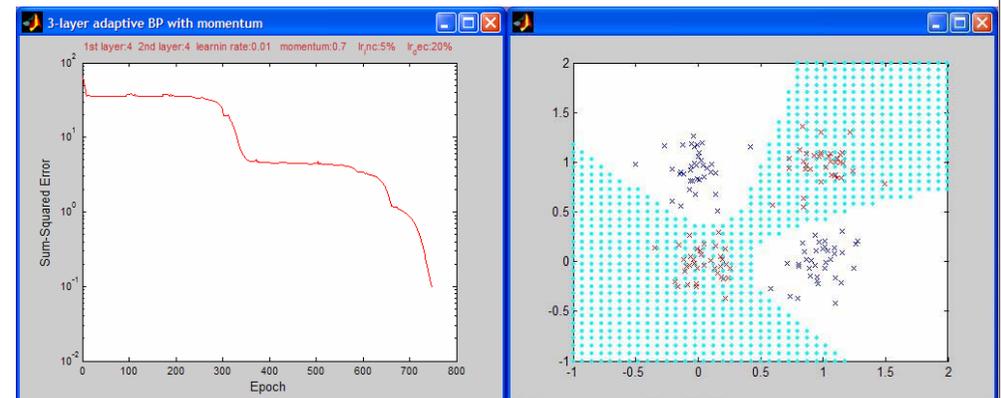
Vereinigung von 2 konvexen Gebieten realisiert mit einem dreischichtigen Netz (2-5-5-1). Mögliche Fehlklassifikation bei Varianz:

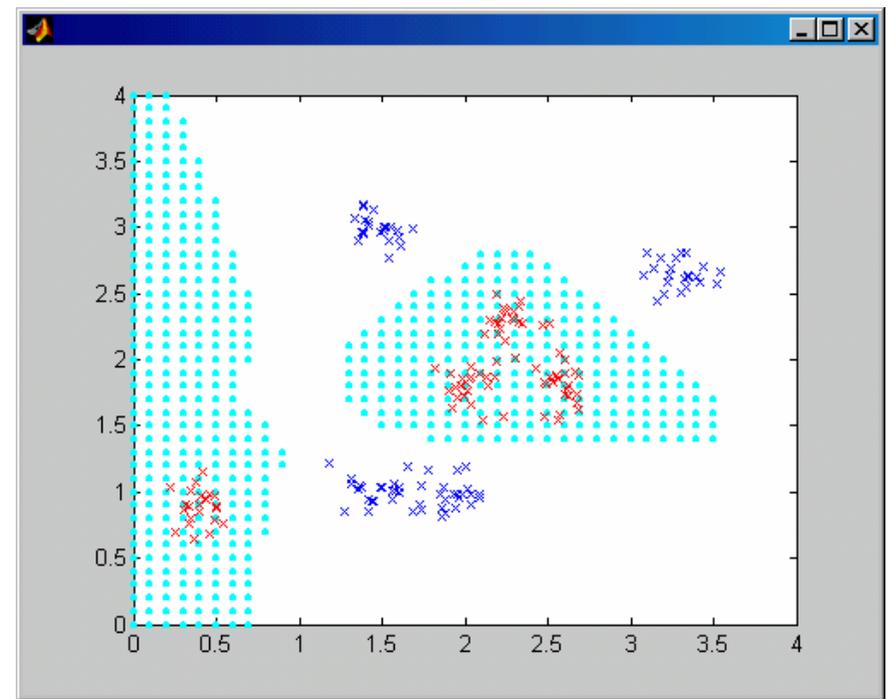
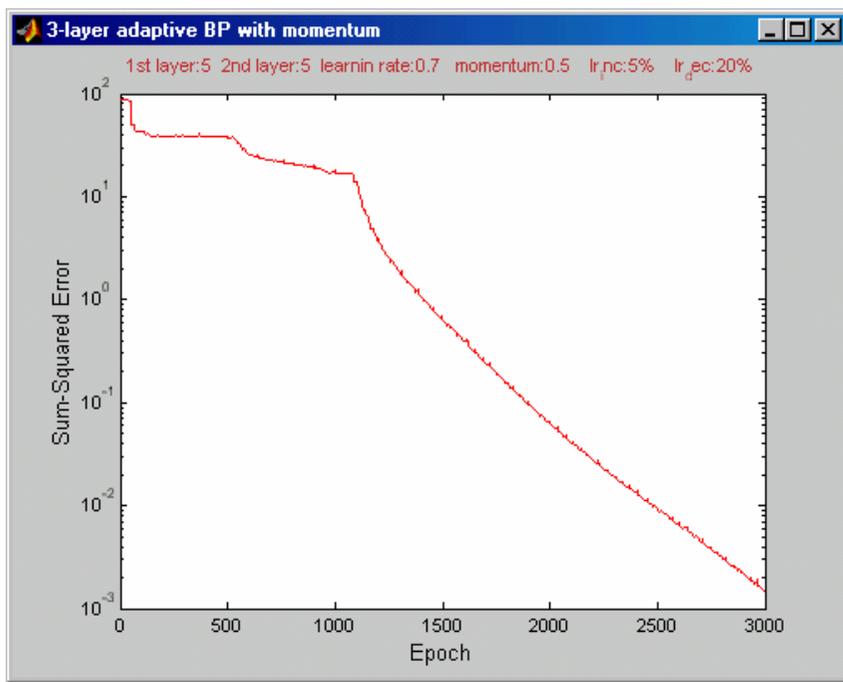
$$\|\mathbf{n}\| > 0,5$$

2-2-2-1 konvergiert leider nicht!



XOR mit hohem Rauschanteil um zweite Lösung zu provozieren!



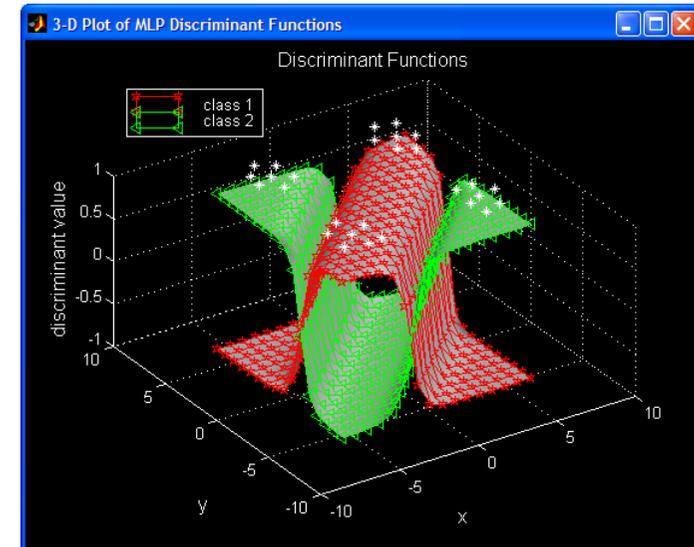


Demo mit MATLAB (Klassifikationgui.m)

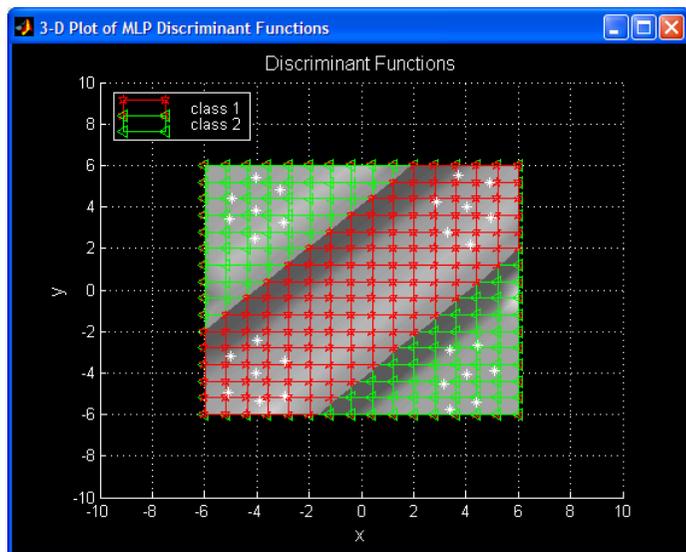
- Öffnen von Matlab
 - zuerst setmypath.m aufrufen, dann
 - C:\Home\ppt\Lehre\ME_2002\matlab\KlassifikatorEntwurf-WinXX\Klassifikationgui
 - Datensatz xor2.mat laden
 - Gewichtsmatrizen xor-trivial.mat laden
 - Gewichtsmatrizen xor-3.mat laden
- Zweiklassenproblem mit Bananenshape eingeben
 - Datensatz Samples/banana_test_samples.mat laden
 - Voreinstellungen aus Models/MLP_5_3_2_banana.mat laden

Start der Matlab-Demo
matlab-Klassifikation_gui.bat

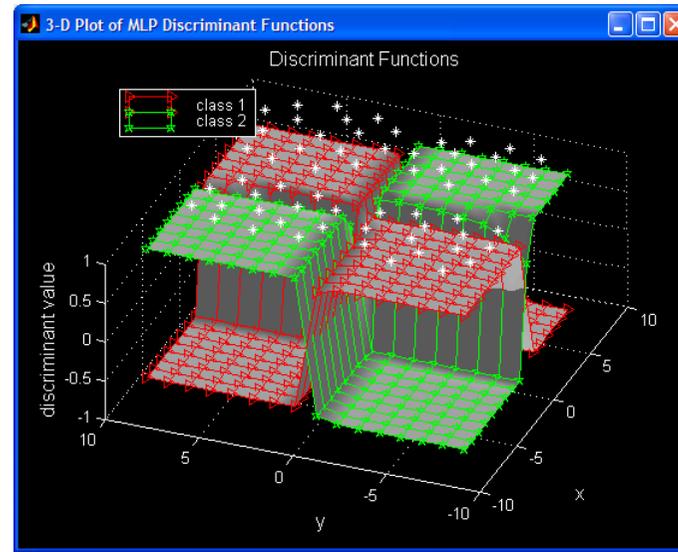
Lösung des XOR-Problems mit zweischichtigem NN



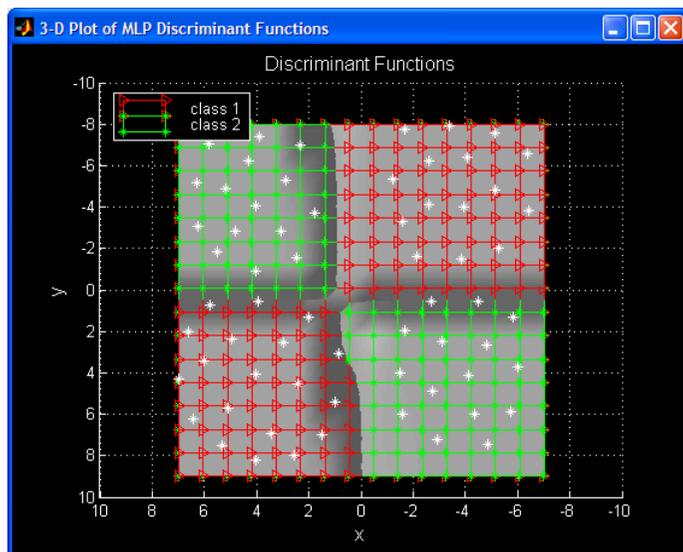
Lösung des XOR-Problems mit zweischichtigem NN



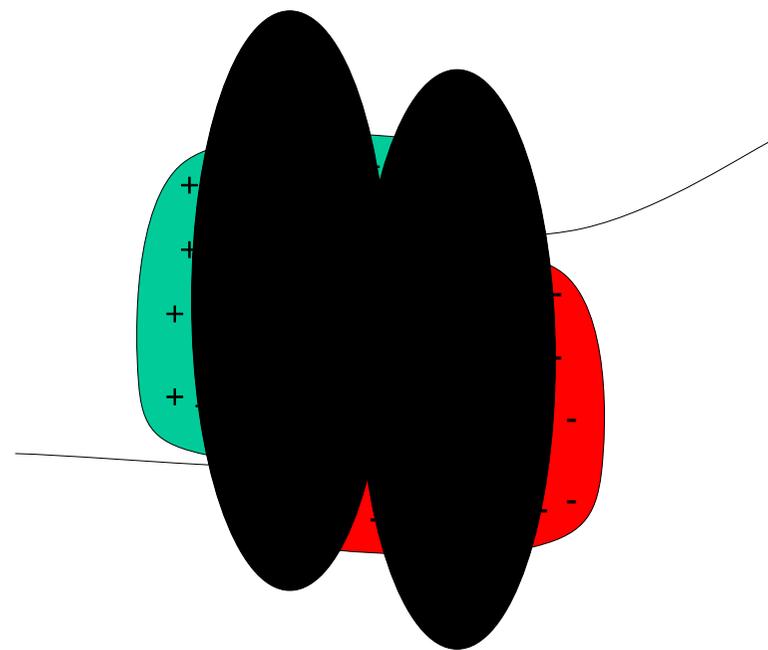
Lösung des XOR-Problems mit dreischichtigem NN

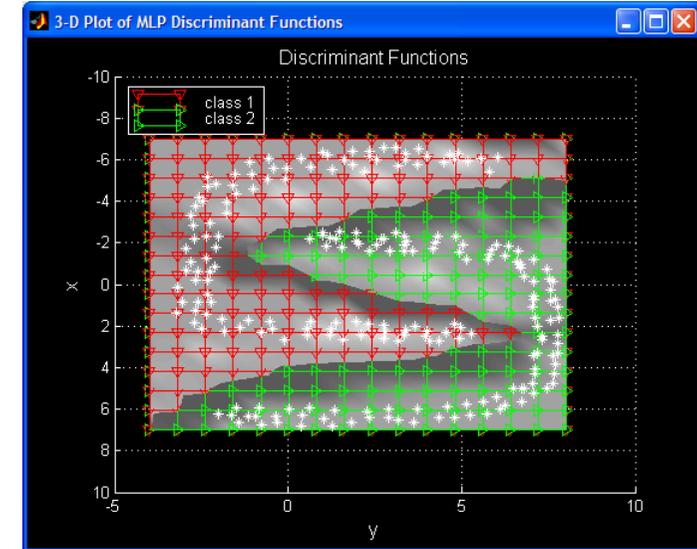
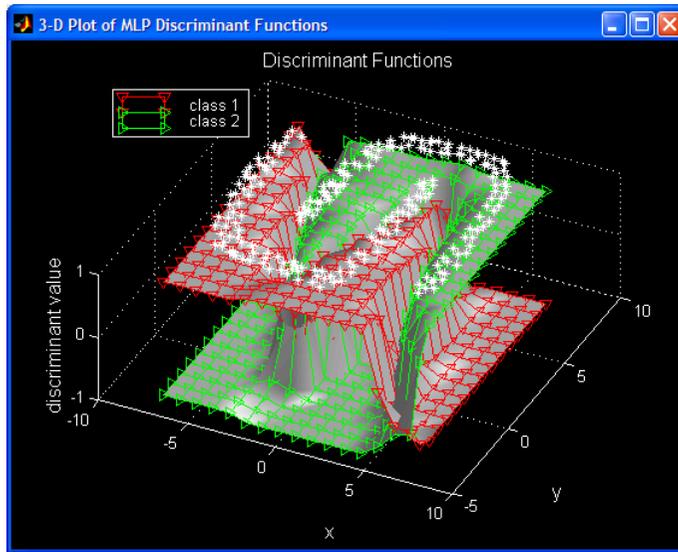


Lösung des XOR-Problems mit dreischichtigem NN



Durch Normalverteilungen schlecht darstellbare Klassen





Konvergenzverhalten von Backpropagation-Algorithmen

- Backpropagation mit gewöhnlichem Gradientenabstieg (steepest descent)

Start der Matlab-Demo
[matlab-BP-gradient.bat](#)

- Backpropagation mit konjugiertem Gradientenabstieg (conjugate gradient)

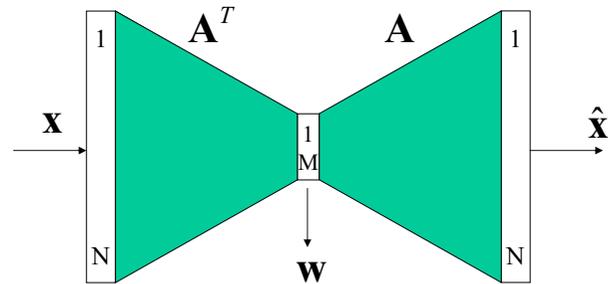
Start der Matlab-Demo
[matlab-BP-CGgradient.bat](#)

NN und deren Eigenschaften

- „Quick and Dirty“. Ein NN-Entwurf ist einfach zu realisieren und man erhält durchaus brauchbare Lösungen (eine suboptimale Lösung ist jedenfalls besser als irgendeine Lösung).
- Es können damit insbesondere auch sehr große Probleme angegangen werden.
- Alle Strategien benutzen jedoch nur „lokale“ Optimierungsfunktionen und erreichen in der Regel kein globales Optimum. Dies ist der gravierende Nachteil für den Einsatz Neuroner Netze.

Verwendung eines NN zur iterativen Berechnung der Hauptkomponentenanalyse (KLT)

- Anstatt die KLT explizit über ein Eigenwertproblem zu lösen, kann auch ein NN zur iterativen Berechnung herangezogen werden.
- Man verwendet ein zweischichtiges Perceptron. Der Ausgang der verdeckten Schicht \mathbf{w} ist der gesuchte Merkmalsvektor.



- Die erste Schicht berechnet den Merkmalsvektor zu:

$$\mathbf{w} = \mathbf{A}^T \mathbf{x}$$

- Dabei wird angenommen, dass eine *lineare Aktivierungsfunktion* zum Ansatz kommt. Die zweite Schicht realisiert die Rekonstruktion von \mathbf{x} mit der transponierten Gewichtsmatrix der ersten Schicht

$$\hat{\mathbf{x}} = \mathbf{A}\mathbf{w}$$

- Das Optimierungsziel ist die Minimierung von:

$$J = E \left\{ \|\hat{\mathbf{x}} - \mathbf{x}\|^2 \right\} = E \left\{ \|\mathbf{A}\mathbf{w} - \mathbf{x}\|^2 \right\} = E \left\{ \|\mathbf{A}\mathbf{A}^T \mathbf{x} - \mathbf{x}\|^2 \right\}$$

- Die folgende Lernregel:

$$\mathbf{A} \leftarrow \mathbf{A} - \alpha (\mathbf{A}\mathbf{w} - \mathbf{x}) \mathbf{w}^T = \mathbf{A} - \alpha \nabla J$$

- führt zu einer Koeffizientenmatrix \mathbf{A} , welche als Produkt von zwei Matrizen geschrieben werden kann (ohne Beweis!):

$$\mathbf{A} = \mathbf{B}_M \mathbf{T}$$

- \mathbf{B}_M ist eine $N \times M$ -Matrix der M dominanten Eigenvektoren von $E\{\mathbf{x}\mathbf{x}^T\}$ und \mathbf{T} eine orthonormale $M \times M$ -Matrix, welche eine Rotation des Koordinatensystems bewirkt, innerhalb eines Raumes, welcher durch die M dominanten Eigenvektoren von $E\{\mathbf{x}\mathbf{x}^T\}$ aufgespannt wird.

- Die Lernstrategie beinhaltet keine Translation. D.h. sie berechnet die Eigenvektoren der Momentenmatrix $E\{\mathbf{x}\mathbf{x}^T\}$ und nicht der Kovarianzmatrix $\mathbf{K} = E\{(\mathbf{x} - \boldsymbol{\mu}_x)(\mathbf{x} - \boldsymbol{\mu}_x)^T\}$. Dies kann jedoch realisiert werden, indem man einen rekursiv geschätzten Erwartungswert $\boldsymbol{\mu}_x = E\{\mathbf{x}\}$ subtrahiert, bevor man die rekursive Schätzung des Merkmalsvektors beginnt. Den gleichen Effekt erzielt man dadurch, dass man einen erweiterten Beobachtungsvektor verwendet:

$$\tilde{\mathbf{x}} = \begin{bmatrix} 1 \\ \mathbf{x} \end{bmatrix} \quad \text{anstatt} \quad \mathbf{x}$$