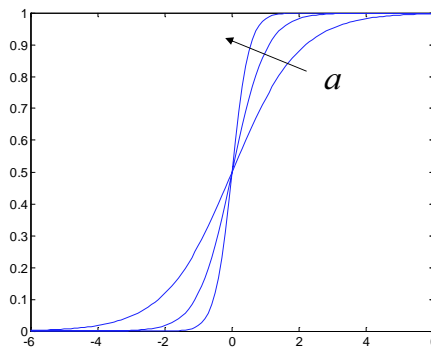


Lernstrategien für Neuronale Netze - Der Backpropagation-Algorithmus

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:

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

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-Perceptron; ohne diese, würde das mehrschichtige in ein triviales lineares einschichtiges Netzwerk zusammenschrumpfen.

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

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

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

Eine Schicht des Neuronalen Netzes

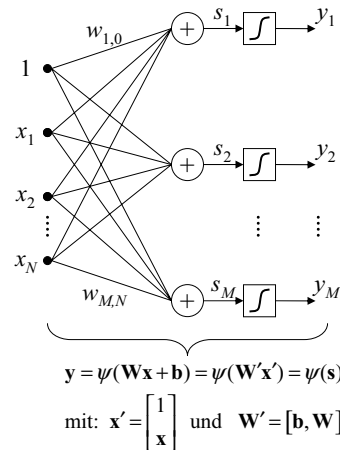
Eine einzige Schicht des NN wird charakterisiert durch die $M \times N$ Koeffizienten der Gewichtsmatrix \mathbf{W} , den Offset-Vektor \mathbf{b} und eine vektorielle Sigmoid-Funktion.

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.

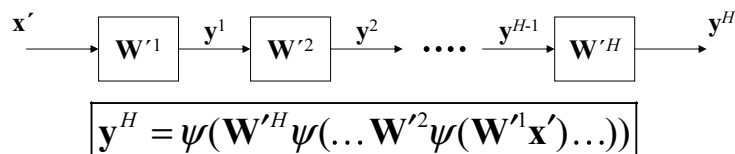


Gestalt der erweiterten Gewichtsmatrix

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

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. Der Erwartungswert ist zu bilden über das Ensemble von $\{\hat{\mathbf{y}}_j, \mathbf{x}_j\}$:

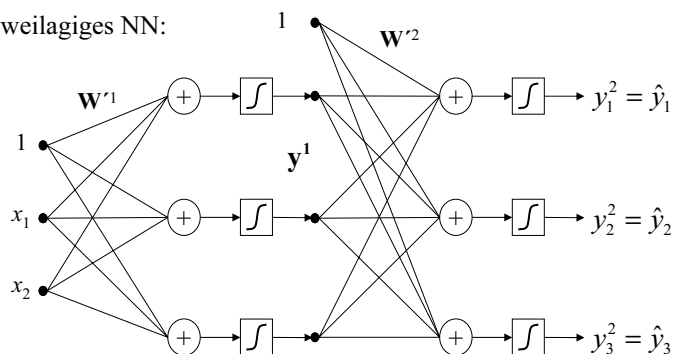
$$J = \min_{\mathbf{w}^i} E \left\{ \|\hat{\mathbf{y}} - \mathbf{y}\|^2 \right\} = \min_{\mathbf{w}^i} E \left\{ \|\mathbf{y}^H - \mathbf{y}\|^2 \right\}$$

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:

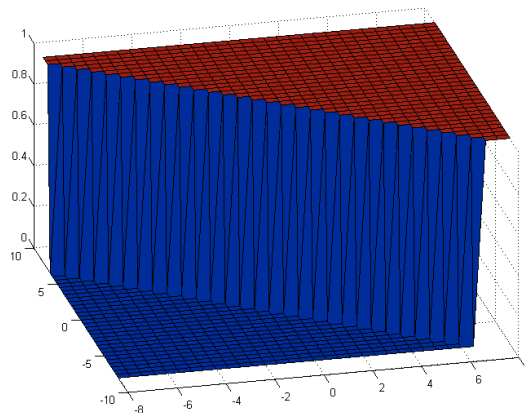


Beziehungen zu dem Konzept der Funktionsapproximation durch eine Linearkombination von Basisfunktionen

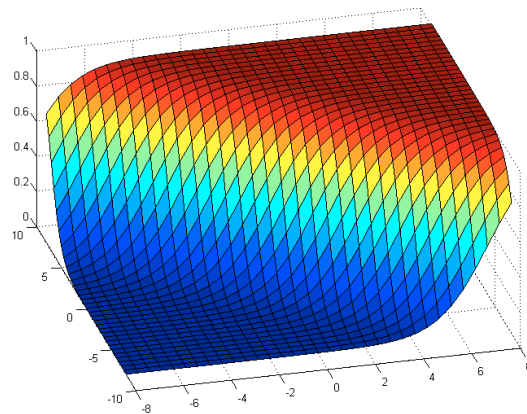
Die erste Schicht erzeugt die Basisfunktionen in Form des verdeckten Vektors \mathbf{y}_1 und die zweite Schicht bildet eine Linearkombination dieser Basisfunktionen.

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

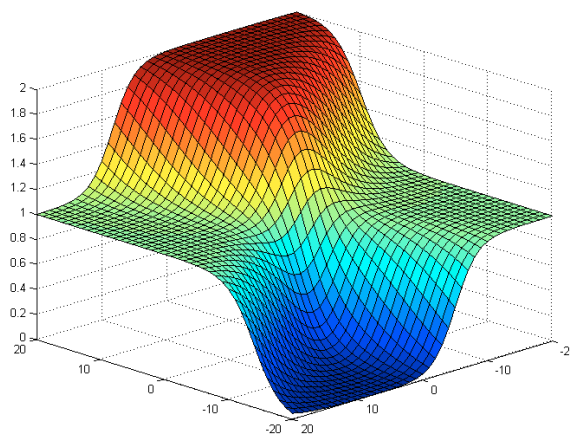
Reaktion eines Neurons mit zwei Eingängen und einer Schwellwertfunktion σ



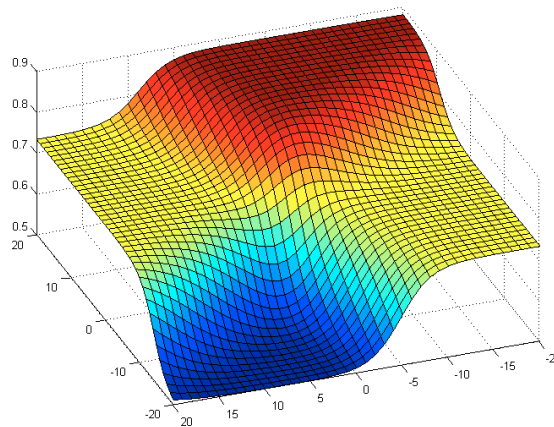
Reaktion eines Neurons mit zwei Eingängen und einer Sigmoidfunktion ψ



Überlagerung zweier Neuronen mit je zwei Eingängen und einer Sigmoidfunktion



Reaktion eines Neurons in der zweiten Schicht auf zwei Neuronen der ersten Schicht jeweils mit Sigmoidfunktionen



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 $\{w'_{ij}{}^h\}$.

- Nach Einsatz der Funktionsabbildung des Multilagenperceptrons erhält man:

$$J = E \left\{ \left\| \underbrace{\psi(\mathbf{W}'^H \dots \psi(\mathbf{W}'^2 \psi(\mathbf{W}'^1 \mathbf{x}')) \dots)}_{\hat{\mathbf{y}}} - \mathbf{y} \right\|^2 \right\}$$

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

$$\boxed{\mathbf{W}' \leftarrow \mathbf{W}' - \alpha \nabla \mathbf{J}} \quad \text{mit: } \mathbf{W}' = \{\mathbf{W}'^1, \mathbf{W}'^2, \dots, \mathbf{W}'^H\}$$

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

- 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

$$\hat{\mathbf{y}}(\mathbf{x}') = \psi(\mathbf{W}'^H \dots \psi(\mathbf{W}'^2 \psi(\mathbf{W}'^1 \mathbf{x}')) \dots)$$

- wird die Kettenregel benötigt.
- Da der Erwartungswert $E\{\dots\}$ nicht zur Verfügung steht, muss er durch den Mittelwert über einzelne Stichproben ersetzt werden:

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

- Der Gradient ergibt sich durch die Mittelung über die Stichprobengradienten:

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

- 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 h-te Schicht:

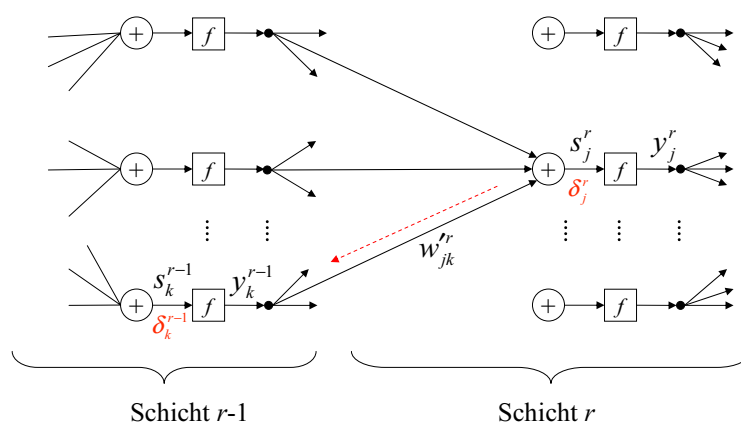
Für die h -te Schicht gilt:

$$y_m^h = \psi(s_m^h) = \psi \left(\sum_{n=0}^{N^h} W_{mn}^h x_n^h \right)$$

x_n^h n-ter Eingang von Schicht h

y_m^h m-ter Ausgang von Schicht h

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:

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

- unter Einführung der *Empfindlichkeit* von Zelle n

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

- ergibt sich:
$$\delta_n^H = -\frac{\partial J}{\partial s_n^H} = -\frac{\partial J}{\partial \hat{y}_n^H} \cdot \frac{\partial \hat{y}_n^H}{\partial s_n^H} = \underbrace{(y_n - \hat{y}_n)}_{\frac{\partial \left(\frac{1}{2} \sum_{k=1}^H (\hat{y}_k - y_k)^2 \right)}{\partial \hat{y}_n}} f'(s_n^H)$$

- und damit schließlich für das Update der Gewichte:

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

- 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

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

- ergibt sich:

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

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

- Man durchläuft alle Lernstichproben, ohne irgendwelche Veränderungen an den Gewichtungskoeffizienten, 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 $[\mathbf{x}'_j, \hat{y}_j]$ sei gegeben
- Berechne die Diskriminierungsfunktion $\hat{\mathbf{y}}$ 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 und berechne daraus rückwärts alle Werte der niederen Schichten gemäß:

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

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

$$f'(x) = \frac{d\psi}{dx} = \psi(x)(1 - \psi(x))$$

- Korrigiere (parallel) alle Perceptron-Gewichte gemäß:

$$w'_{nm} \leftarrow w'_{nm} - \alpha \frac{\partial J}{\partial w'_{nm}} = w'_{nm} - \alpha \delta_n^r y_m^{r-1} \quad \begin{array}{l} h = 1, 2, \dots, H \\ \text{für } n = 1, 2, \dots, N^h \\ m = 1, 2, \dots, M^h \end{array}$$

- Beim individuellen Lernen werden die Gewichte $\{w'_{nm}{}^h\}$ 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:

$$w'_{nm} \leftarrow w'_{nm} - \sum_i \alpha \delta_n^r(i) y_m^{r-1}(i) \quad \begin{array}{l} h = 1, 2, \dots, H \\ \text{für } n = 1, 2, \dots, N^h \\ m = 1, 2, \dots, M^h \end{array}$$

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:

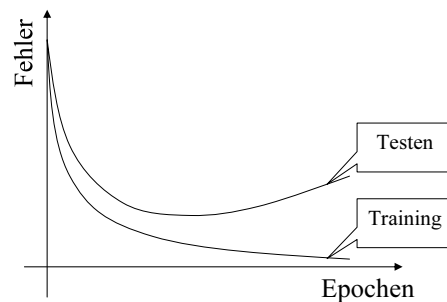
$$T = \dim(\mathbf{W}) = \sum_{h=1}^H (N^{h-1} + 1) N^h \quad \text{mit } N^0 = N \quad \text{und } N^H = K$$

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

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

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

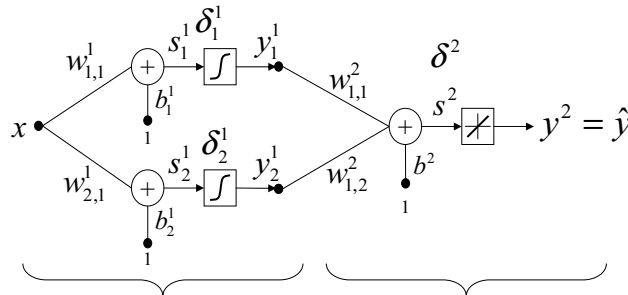
- Wenn $T = \dim(\mathbf{W})$ die Anzahl der Gewichte und Biasterme ist, so läßt sich einfach nachvollziehen, daß $O(T)$ Berechnungsschritte für die Vorwärtssimulation benötigt werden, $O(T)$ für das Backpropagieren des Fehlers und ebenso $O(T)$ Operationen für die Korrektur der Gewichte, also erhält 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:

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

- 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“)



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

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

- 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 mit linearer Aktivierungsfunktion ergibt sich:

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

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

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

- 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:

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

nnd11bc
File Edit View Insert Tools Window Help

Neural Network DESIGN Backpropagation Calculation

Last Error: ???

Input: p = 1.0

Target: t = 1+sin(p*pi/4) = 1.707

Simulate:
a1 = logsig(W1*p+b1) = [0.321; 0.368]
a2 = purelin(W2*a1+b2) = 0.446
e = t-a2 = 1.261

Backpropagate:
s2 = 2*dpurelin(n2)/dn2*e = -2.522
s1 = dlogsig(n1)/dn1*W2*s2 = [-0.049; 0.100]

Update:
W1 = W1-lr*s1*p' = [-0.265; -0.420]
b1 = b1-lr*s1 = [-0.475; -0.140]
W2 = W2-lr*s2*a1' = [0.171; -0.077]
b2 = b2-lr*s2 = 0.732

Continue

Reset

Random

Contents

Close

Chapter 11

H. Burkhardt, Institut für Informatik, Universität Freiburg
ME-II, Kap. 8b 27

Backpropagation zum Trainieren eines zweischichtigen Netzwerks zur Funktionsapproximation

(Matlab-Demo: „Function Approximation“)

- Gesucht wird eine Approximation der Funktion

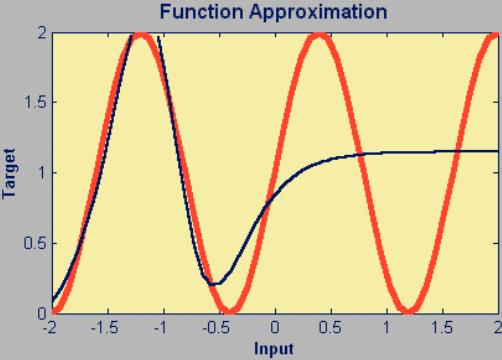
$$y(x) = 1 + \sin\left(i \cdot \frac{\pi}{4} x\right) \quad \text{für} \quad -2 \leq x \leq 2$$

- Mit zunehmenden Wert von i steigen die Anforderungen an das MLPC-Netzwerk. Die auf der Sigmoid-Funktion aufbauende Approximation benötigt mehr und mehr Schichten, um mehrere Perioden der Sinus-Funktion darzustellen.

nnd11fa

File Edit View Insert Tools Window Help

Neural Network DESIGN **Function Approximation**



Click the [Train] button to train the logsig-linear network on the function at left.

Use the slide bars to choose the number of neurons in the hidden layer and the difficulty of the function.

Train

Contents

Close

Chapter 11

Number of Hidden Neurons S1: 2

1 9

Difficulty Index: 5

1 9