

Lernen und Klassifizieren

AS2-2

Assoziativspeicher

Lineare Klassifikation

Lernen linearer Klassifikation

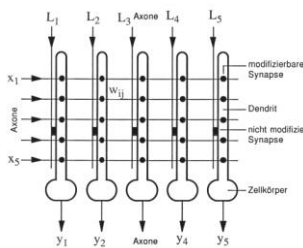
Lernen und Zielfunktion

Stochast. Klassifikation

Lernen in Multilayer-Netzen

Backpropagation-Lernen

Neuro-Modell des Assoziativspeichers



Funktion:

Jede Komp.ist lin. Summe

$$z_i = w_i x$$

Nichtlin. Ausgabe:

$$y_i = S_B(z_i) = \begin{cases} 1 & z > \theta \\ 0 & z \leq \theta \end{cases}$$

Lernen von W ?

Lernen im Assoziativspeicher

- Speichern aller N Muster mit Hebbischer Regel

$$w_{ij} = \sum_{k=1}^N \Delta w_{ij} = \sum_{k=1}^N \gamma_k L_i^k x_j^k \quad w_{ij}(0) := 0$$

- Auslesen eines Musters r

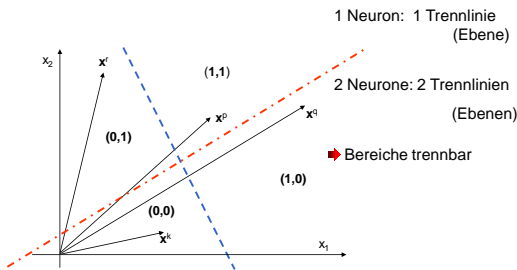
$$y = \mathbf{W}\mathbf{x}^r = \mathbf{z} = \gamma_r \mathbf{L}^r(\mathbf{x}^r)^T \mathbf{x}^r + \sum_{k \neq r} \gamma_k \mathbf{L}^k (\mathbf{x}^k)^T \mathbf{x}^r$$

assoziierte Antwort + Übersprechen von anderen Mustern

- Orthogonale Muster \mathbf{x}^r :** Übersprechen = 0, exakte Reproduktion.
- Nicht-orthogonale Muster:** Schwellwerte nötig zum Unterdrücken des Übersprechens.

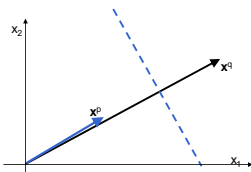
Trennung mehrerer Klassen

Erinnerung: Lineare Separierung



Trennung mehrerer Klassen

Problem: Klassenentscheidung über Korrelationsgröße



Entscheidung über x :

Klasse p: $\mathbf{x}\mathbf{x}^p > \mathbf{x}\mathbf{x}^q$

Klasse q: $\mathbf{x}\mathbf{x}^p < \mathbf{x}\mathbf{x}^q$

Frage: $x = \mathbf{x}^p$:
In welche Klasse?

Antwort:
in Klasse q !

Lösung

$(\mathbf{x}-\mathbf{y})^2 = \mathbf{x}^2 - 2\mathbf{x}\mathbf{y} + \mathbf{y}^2$ ist minimal $\Leftrightarrow \mathbf{x}\mathbf{y}$ ist maximal
genau dann, wenn

Konstante Länge $c = |\mathbf{x}|=|\mathbf{y}|$ (normierte Musteraktivität)

Trennung mehrerer Klassen

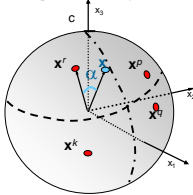
Erweiterung der Mustertupel

$$\mathbf{x} \rightarrow \mathbf{X}' = (x_0, x_1, x_2, \dots, x_n) \text{ mit } |\mathbf{x}'| = \text{const}$$

$$\text{weil } x_0^2 = c^2 - |(x_1, x_2, \dots, x_n)|^2 \geq 0 (!)$$

≅ Einbettung in den Hyperraum

Beispiel: 2-dim → 3-dim



Entscheidung durch

$$\cos(\alpha) = \frac{\mathbf{x}^i \mathbf{x}^f}{|\mathbf{x}^i| |\mathbf{x}^f|} = c^{-2} \mathbf{x}^i \mathbf{x}^f$$

$\cos(\alpha)$ monoton fallend \Rightarrow
Winkel als Distanzmaß

$\min \alpha \Leftrightarrow \max$ Korrelation

Assoziativspeicher: Speicherkapazität

M Tupel (\mathbf{x}, \mathbf{y}) gegeben:

Wie viele können zuverlässig gespeichert werden?

$\mathbf{x}^1 = \mathbf{x}^2 = \dots = \mathbf{x}^M$: nur ein Muster speicherbar.

$\mathbf{y}^1 = \mathbf{y}^2 = \dots = \mathbf{y}^M$: beliebig viele Muster speicherbar, da Antwort \mathbf{y} immer richtig.

Problem der **Kodierung** der Muster ! Sei $|\mathbf{x}| = a$.

- Maximaler Musterabstand
 $\max d(\mathbf{x}^p, \mathbf{x}^q) = \min \mathbf{x}^p \mathbf{x}^q = 0$ bei orthogonalen Mustern
Reelle Komponenten: n Dimensionen $\Rightarrow n$ orthogonale Basisvektoren
Binäre Komponenten:

$$M_{\max} = \left\lfloor \frac{n}{a} \right\rfloor \quad \text{z.B. } n=100, a=10, \text{ also max } M=10$$

- Mittlere Abstand maximal

$$\frac{\sum d(\mathbf{x}^i, \mathbf{x}^j)}{\binom{M}{2}} \Big|_{M=a} = 0 \Rightarrow M = \binom{n}{a} = \binom{n}{n/2} \quad \text{z.B. } n=100$$

$$\text{max } M \approx 2^{n/a} \approx 2^{10} \approx 10^{29}$$

Assoziativspeicher: Binärspeicher

Spärliche Kodierung Binäre Muster

Speichern: $w_{ij} = \bigvee_p y_j^p x_i^p = \max_p y_j^p x_i^p$

Kapazität: $H_B = \ln 2 = 0,693$ Bit pro Speicherzelle
vergleichbar mit CAM-Speicher

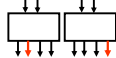
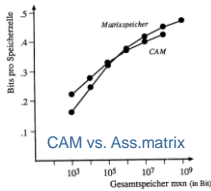
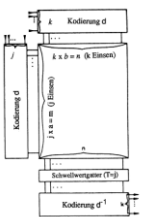
Palm 1980

Kodierung

$$k = a_x = \text{ld } m$$

$$j = a_y = O(\log n)$$

Konstante Zahl
von 1 durch eine
Leitung pro
Eingabecode



Assoziativspeicher

Lineare Klassifikation

Lernen linearer Klassifikation

Lernen und Zielfunktion

Stochast. Klassifikation

Lernen in Multilayer-Netzen

Backpropagation-Lernen

Rüdiger Brause: Adaptive Systeme AS-2 WS 2013

Klassenbildung heute

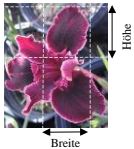
Objekte werden durch **Merkmale** beschrieben

- **z.B. qualitativ** Mensch = (groß, braune Augen, dunkle Haare, nett, ...)
- **quantitativ** Mensch = (Größe=1,80m, Augenfarbe=2, Haarfarbe=7, ...)

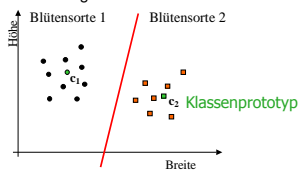
Idee = Form = „Klassenprototyp“

Muster eines Objekts

≡ (Breite, Höhe) = x



Trennung von **Klassen**



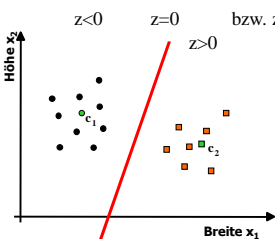
Klassifizierung = Ermitteln der Geradengleichung bzw Parameter c_1, c_2 .

Rüdiger Brause: Adaptive Systeme AS-2 WS 2013

- 11 -

Klassentrennung

Klassentrennung durch Trenngerade mit $f(x_1) = x_2 = w_1x_1 + w_2$



bzw. $z = w_1x_1 + w_2x_2 + w_3x_3 = 0$
mit $x_3 := 1$

$$\text{Mit } z = \sum_{i=1}^n w_i x_i = \mathbf{w}^T \mathbf{x}$$

Klassenentscheidung

$$y = S(z) = \begin{cases} 0 & z \leq 0 & \mathbf{x} \text{ aus Klasse 1} \\ 1 & z > 0 & \mathbf{x} \text{ aus Klasse 2} \end{cases}$$

Rüdiger Brause: Adaptive Systeme AS-2 WS 2013

Trennung mehrerer Klassen

•DEF Lineare Separierung

Seien Muster \mathbf{x} und Parameter \mathbf{w} gegeben.

Zwei Klassen Ω_1 und Ω_2
des Musterraums $\Omega = \Omega_1 \cup \Omega_2$ mit $\Omega_1 \cap \Omega_2 = \emptyset$

heißen **linear separierbar**,

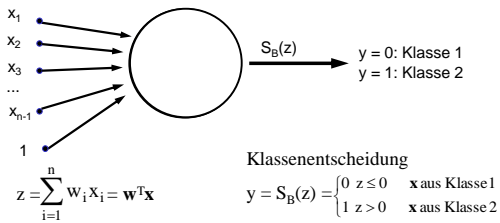
falls eine Hyperebene $\{\mathbf{x}^*\}$ existiert mit $g(\mathbf{x}^*) = \mathbf{w}^T \mathbf{x}^* = 0$,

so daß für alle $\mathbf{x} \in \Omega_1$ gilt $g(\mathbf{x}) < 0$

und für alle $\mathbf{x} \in \Omega_2$ gilt $g(\mathbf{x}) > 0$.

Klassentrennung durch formales Neuron

Klassentrennung durch binäres Neuron



ALSO :

WIE erhalten wir die richtigen
Gewichte,

d.h. die richtige Klassifizierung

?

Lernen !

Assoziativspeicher

Lineare Klassifikation

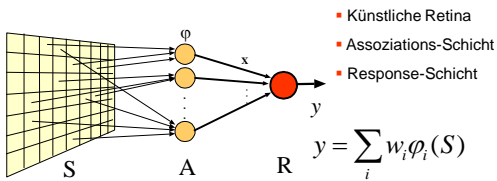
Lernen linearer Klassifikation

- Lernen und Zielfunktion
- Stochast. Klassifikation
- Lernen in Multilayer-Netzen
- Backpropagation-Lernen

Rüdiger Brause: Adaptive Systeme AS-2 WS 2013

Das Perzeptron

Idee: Reize wiedererkennen
Rosenblatt 1958



- Verbindungen zu A fix (zufällig): $\mathbf{x} = (x_1, \dots, x_n)^T = (\phi_1(S), \dots, \phi_n(S))^T$
- Stärke der Verbindungen zu R veränderbar: $\mathbf{w} = (w_1, \dots, w_n)^T$

Rüdiger Brause: Adaptive Systeme AS-2 WS 2013

Das Perzeptron

Entscheiden $\Omega := \{\mathbf{x}\}$ alle Muster, $\Omega = \Omega_1 + \Omega_2$
 Ω_1 : Menge aller \mathbf{x} aus Klasse 1
 Ω_2 : Menge aller \mathbf{x} aus Klasse 2

DEF Log. Prädikat $y = \begin{cases} \text{TRUE} & \text{wenn } \mathbf{w}^T \mathbf{x} > s \\ \text{FALSE} & \text{wenn } \mathbf{w}^T \mathbf{x} \leq s \end{cases}$ Schwelle

Mit den Erweiterungen $\mathbf{x} = (x_1, \dots, x_n, 1)^T$
 $\mathbf{w} = (w_1, \dots, w_n, s)^T$

wird $y = \begin{cases} \text{TRUE} & \text{wenn } \mathbf{w}^T \mathbf{x} > 0 \\ \text{FALSE} & \text{wenn } \mathbf{w}^T \mathbf{x} \leq 0 \end{cases}$

Rüdiger Brause: Adaptive Systeme AS-2 WS 2013

Das Perzeptron: Pseudo-code 3

DEF numerische Werte
$$y = \begin{cases} 1 & \text{wenn } \mathbf{w}^\top \mathbf{x} > 0 \\ 0 & \text{wenn } \mathbf{w}^\top \mathbf{x} \leq 0 \end{cases}$$

$$L(\mathbf{x}) = \begin{cases} 1 & \text{wenn } \mathbf{x} \in \Omega_1 \\ 0 & \text{sonst} \end{cases}$$

PERCEPT3:

Wähle zufällige Gewichte \mathbf{w} zum Zeitpunkt $t=0$.

REPEAT

$t = t + 1$;

$$\mathbf{w}(t) = \mathbf{w}(t-1) + \gamma (L(\mathbf{x}) - y(\mathbf{x})) \mathbf{x}(t) \quad \text{Fehler-Lernregel}$$

UNTIL (alle \mathbf{x} richtig klassifiziert)

Sogar ohne Umdefinition der Muster aus Ω_2 !

Das Perzeptron: Konvergenz

Perzeptron - Konvergenztheorem (Minsky Papert 1988)

Wenn die Mustermenge Ω_1 **linear separierbar** ist, so konvergiert der Algorithmus bei $t \rightarrow \infty$

Problem: Wenn Klassen sich überlappen, so wird die Grenzlinie bei $\gamma = 1$ immer hin und her geschoben

Das Perzeptron: Zielfunktion

Ziel: Verallgemeinerung der Lernregel

Hier: Minimierung **aller** Fehlentscheidungen

DEF
$$R(\mathbf{w}) = \sum_{\mathbf{x} \in \Omega_F} -\mathbf{w}^\top \mathbf{x} \quad \text{Perzeptron-Zielfunktion}$$

„Energie“

Neuformulierung erwartetes Lernen:

$$\mathbf{w}(t) = \mathbf{w}(0) - \gamma \nabla_{\mathbf{w}} R(\mathbf{w}) \quad \text{Gradient}$$

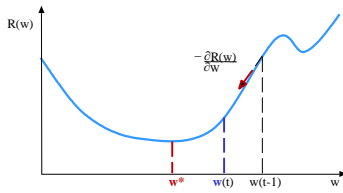
d.h.

$$\mathbf{w}(t) = \mathbf{w}(0) + \gamma \sum_{\mathbf{x} \in \Omega_F} \mathbf{x}$$

$$\mathbf{w}(t) = \mathbf{w}(t-1) - \gamma \nabla_{\mathbf{w}_{t-1}} r(\mathbf{w}(t-1), \mathbf{x}) \quad \text{Stochast. Lernen}$$

Lernen durch Iteration

Gradientenabstieg einer Zielfunktion $R(w)$



$$\Delta \mathbf{w} := (\mathbf{w}(t-1) - \mathbf{w}(t)) \sim -\nabla_{\mathbf{w}} R(\mathbf{w}(t-1))$$

$$\mathbf{w}(t) = \mathbf{w}(t-1) - \gamma(t) \nabla_{\mathbf{w}} R(\mathbf{w}(t-1))$$

Was kann ein Perzeptron ?

Erwartung: „Intelligente Leistungen“ durch Wahl von $\phi(S)$

Abbildung der Merkmale auf linear separierbare
Muster Mengen

Perzeptronarten

- **diameter-limited perceptrons**
nur Bildpunkte aus einem begrenzten Radius
- **order-restricted perceptrons**
von maximal n (beliebigen) Bildpunkten abhängig
- **random perceptrons**
eine zufällige Auswahl aller Bildpunkte

Was kann ein Perzeptron ?

Topologische Prädikate, z.B.

- „X ist ein Kreis“ ?
- „X ist eine konvexe Figur“ ?
- „X ist eine zusammenhängende Figur“ ?
- ...

Tatsache: *keine* korrekte Klassifizierung von
Punktmengen X (Bildpixeln) dieser Arten

Nur „X hat Eulerzahl E“

$$E(X) = K(X) - \text{Anzahl der Löcher}$$

Was kann ein Perzeptron ?

Eulerzahl E

$E(X) := K(X) - \text{Anzahl der Löcher}$

$K(X) := \text{zusammenhängende Komponenten}$

Loch := zusamm. Komponente der komplementären Menge



$$K(x) = 2, \text{ Löcher} = 1 \Rightarrow E(x) = 1$$

Was kann ein Perzeptron ?

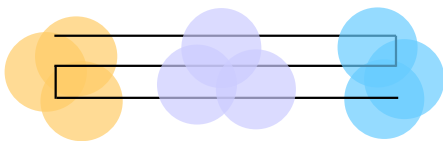
Beispiel: keine korrekte Klassifizierung von Punktfolgen X (Bildpixeln) für Prädikat „ X ist Typ A“ möglich mit „diameter-limited“ Perzeptron

		Typ A
Muster 1	Muster 2	
		Nicht Typ A
Muster 3	Muster 4	

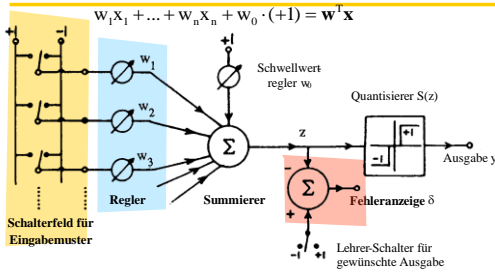
Was kann ein Perzeptron ?

Beweis:

offen: Typ A	Nicht Typ A



Adaline: Aktivität



$$w_1 x_1 + \dots + w_n x_n + w_0 \cdot (+1) = \mathbf{w}^T \mathbf{x}$$

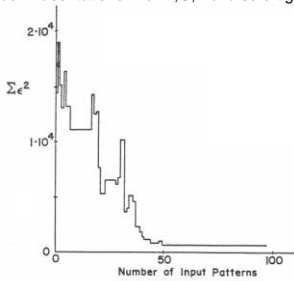
$$w_i \in \mathbb{R}^+$$

$$x_j \in \{-1, +1\}$$

$$y = \begin{cases} +1 & \text{wenn } \mathbf{w}^T \mathbf{x} > 0 \\ -1 & \text{wenn } \mathbf{w}^T \mathbf{x} \leq 0 \end{cases}$$

Adaline: Aktivität

Verlauf des Klassifizierungsfehlers für „Klasse T liegt vor“ bei Präsentationen von T,G,F und sofortiger Nachregelung



Adaline: Lernalgorithmus

Minimierung des erwarteten quadratischen Fehlers

$$R(\mathbf{w}, L) := \langle (z(\mathbf{x}) - L(\mathbf{x}))^2 \rangle_{\mathbf{x}} = \langle (\mathbf{w}^T \mathbf{x} - L(\mathbf{x}))^2 \rangle_{\mathbf{x}}$$

durch Anpassung der Parameter

$$\mathbf{w}(t) = \mathbf{w}(t-1) - \gamma(t) \frac{\partial R(\mathbf{w}(t-1))}{\partial \mathbf{w}}$$

$$\mathbf{w}(t) = \mathbf{w}(t-1) - \gamma(t)(\mathbf{w}^T \mathbf{x} - L(\mathbf{x}))\mathbf{x} \quad \text{stochastische Approximation}$$

$$\mathbf{w}(t) = \mathbf{w}(t-1) - \gamma(t) \frac{\mathbf{x}}{|\mathbf{x}|^2} (\mathbf{w}^T \mathbf{x} - L(\mathbf{x})) \quad \text{Widrow-Hoff Lernregel}$$

Übersicht: Lernen

Assoziativspeicher

1. Muster \mathbf{x}^k eingespeichert

$$\mathbf{w}_i(1) = L_i^k \mathbf{x}^k \quad (\text{Hebb'sche Regel})$$

Perzeptron

$$\mathbf{w}_i(t) = \mathbf{w}_i(t-1) + \gamma(L_i(\mathbf{x}) - y_i)\mathbf{x} \quad (\text{Fehler-Lernregel})$$

$$\mathbf{w}_i(1) = (L_i(\mathbf{x}^k) - y_i)\mathbf{x}^k = L_i^k \mathbf{x}^k \quad \text{bei } \mathbf{w}_i(0) = 0 \Rightarrow y_i^k(0) = 0.$$

Adaline

$$\mathbf{w}_i(t) = \mathbf{w}_i(t-1) + \gamma(t)(L(\mathbf{x}) - z_i)\mathbf{x} \quad (\text{Gradientenabstieg})$$

$$\mathbf{w}_i(1) = (L_i(\mathbf{x}^k) - z_i)\mathbf{x}^k = L_i^k \mathbf{x}^k \quad \text{bei } \mathbf{w}_i(0) = 0 \Rightarrow z_i^k(0) = 0.$$

Assoziativspeicher = Grundfunktion von Netzen

Assoziativspeicher

Lineare Klassifikation

Lernen linearer Klassifikation

Lernen und Zielfunktion

Stochast. Klassifikation

Lernen in Multilayer-Netzen

Backpropagation-Lernen

Rüdiger Brause: Adaptive Systeme AS-2 WS 2013

Übersicht Lernarten

- **Beispiel-basiertes Lernen** (*example based learning, feedback learning*)
Gegeben: (Eingabe \mathbf{x} , gewünschte Ausgabe L)
Ziel: Differenz zwischen y und L im Laufe des Lernens klein machen.
- **Erklärungs-basiertes Lernen** (*explanation based learning EBL*)
Gegeben: **Beispielpaare, Ziel** sowie **Regeln**, es zu erreichen.
Lernen: Generalisierung der Beispiele.
(regelbasierte Systeme, nicht bei neuronalen Netzen)
- **Score-basiertes Lernen** (*reinforcement learning*)
Gegeben: **skalares Gütemaß** ("gut", "schlecht", mit Abstufungen dazwischen) für Lernleistung.
Lernen: ?? *Der Lernende muss daraus selbst sehen, was an der Ausgabe zu ändern ist.*
- **Unüberwachtes Lernen** (*observation based learning, emotion based learning, similarity learning*)
Gegeben: **keine explizite** Rückmeldung über die Güte seines Lernens
Lernen: Vergleich gewünschte Auswirkungen mit beobachteten Auswirkungen. Folgerung für geeignete Verhaltensänderung.

Lernen durch Iteration

Modifikationen Gradientenabstieg

Taylorentwicklung

$$f(x+\Delta x) = f(x) + \frac{\partial}{\partial x} f(x) \Delta x + \frac{1}{2!} \frac{\partial^2}{\partial x^2} f(x) (\Delta x)^2 + \dots$$

$$R(\mathbf{w}+\Delta \mathbf{w}) - R(\mathbf{w}) = \nabla_{\mathbf{w}} R(\mathbf{w})^T \Delta \mathbf{w} + \frac{1}{2} \Delta \mathbf{w}^T \mathbf{R} \Delta \mathbf{w} + \dots$$

mit $\mathbf{R} = \frac{\partial^2 R}{\partial w_i \partial w_j}$ *Hesse-Matrix*

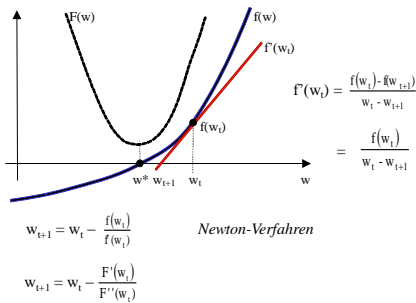
Conjugate gradient

$$R(\mathbf{w}+\Delta \mathbf{w}) - R(\mathbf{w}) = (\nabla_{\mathbf{w}} R(\mathbf{w}))^T + \frac{1}{2} \Delta \mathbf{w}^T \mathbf{R} \Delta \mathbf{w} = 0$$

löse n -dim Gleichungssystem für $\Delta \mathbf{w}$

Lernen durch Iteration

Newton-Iteration



Lernen durch Iteration

Konvergenz des Gradientenverfahrens

Es ist $R(t) =$ *Ljapunov-Funktion* mit Konvergenz, wenn

- $R(t+1) \leq R(t)$ bzw. $\frac{\partial R}{\partial t} \leq 0$ *monoton fallend*
- Ex. endliches $R_0 \leq R(t)$ für jedes t *Ljapunov-Bedingung*

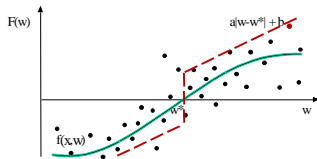
Also: **Wenn** $\frac{\partial R}{\partial t}(w(t)) = \nabla_{\mathbf{w}} R(\mathbf{w}) \frac{\partial \mathbf{w}}{\partial t} \leq 0$ **dann Konvergenz**

Hinreichend dafür: $\frac{\partial \mathbf{w}}{\partial t} = -\gamma \nabla_{\mathbf{w}} R(\mathbf{w})$ mit $\gamma > 0$
weil $\frac{\partial R}{\partial t} = -\gamma (\nabla_{\mathbf{w}} R(\mathbf{w}))^2 \leq 0$

Mit $\frac{\partial \mathbf{w}}{\partial t} \approx \frac{\Delta \mathbf{w}}{\Delta t}$ und $\Delta t = 1$ ist
 $\mathbf{w}(t) - \mathbf{w}(t-1) = -\gamma \nabla_{\mathbf{w}} R(\mathbf{w})$ *Gradientenabstieg*

Stochastische Approximation

Gesucht: Nullstelle einer stochast. Funktion $f(x,w) = R'(x,w)$



Methode 1: Alle Ereignisse x abwarten und dann $F(w) = \langle f(x,w) \rangle_x$ bilden
 $w(t) = w(t-1) - \gamma(t) F(w(t-1))$

Methode 2: Einfach $f(x,w)$ verwenden *Robbins, Monro 1951*
 $w(t) = w(t-1) - \gamma(t) f(w(t-1), x(t))$ *stochastische Approximation*

Stochastisches Lernen

Lernen mit Zielfunktion $R(w) = \langle r(w,x) \rangle_x$

$$w(t) = w(t-1) - \gamma(t) \nabla_w R(w(t-1))$$

wird ersetzt durch

Lernen mit stochast. Zielfunktion $r(w,x)$

$$w(t) = w(t-1) - \gamma(t) \nabla_w r(w(t-1), x(t))$$
 stochastisches Lernen

Stochastische Approximation

Voraussetzungen *das klein Gedruckte...*

- die Funktion $F(w) := \langle f(x,w) \rangle_x$ ist *zentriert*, d.h. $F(w^*) = 0$
- $F(w)$ ist *ansteigend*, d.h. $F(w < w^*) < 0, F(w > w^*) > 0$.
- $F(w)$ ist *beschränkt* mit $|F(w)| \leq a|w-w^*|+b < \infty$ $a, b > 0$
- $f(x,w)$ hat *endliche Varianz*, d.h. $\sigma^2(w) = \langle (F(w) - f(x,w))^2 \rangle_x < \infty$
- $\gamma(t)$ *verschwindet*, $\gamma(t) \rightarrow 0$
- $\gamma(t)$ wird nicht zu schnell *klein* $\sum_{t=1}^{\infty} \gamma(t) = \infty$
- $\gamma(t)$ wird nicht zu *groß* $\sum_{t=1}^{\infty} \gamma(t)^2 < \infty$

Dann ex.

$$\lim_{t \rightarrow \infty} \langle (w(t) - w^*)^2 \rangle = 0 \quad \text{mittl. quadr. Konv.} \quad \text{Robbins-Monro}$$

$$P(\lim_{t \rightarrow \infty} w(t) = w^*) = 1 \quad \text{Blum}$$

Stochastische Iteration: Konvergenz

Beispiel

Sei die Zufallsvariable X gegeben, geschätzt durch W .

Abweichung bei der Schätzung ist

$$R(w) = \langle r(w, X) \rangle_x = \langle (w - X)^2 \rangle_x \quad \text{mean squared error}$$

$$\mathbf{w}(t) = \mathbf{w}(t-1) - \gamma(t) \nabla_{\mathbf{w}} r(\mathbf{w}(t-1), \mathbf{X}(t)) \quad \text{stoch. Gradient}$$

$$\mathbf{w}(t) = \mathbf{w}(t-1) - \gamma(t)(\mathbf{w}(t-1) - \mathbf{X}(t))$$

Zeitabhängigkeit

$R(w) \rightarrow R(w^*)$ bei $w \rightarrow w^*$

stoch. und erwarteter Verlauf?

Stochastische Iteration: Konvergenz

Stochastische Iteration

$$\mathbf{w}_i(t) = \mathbf{w}_i(t-1) - \gamma(t)(\mathbf{w}_i(t-1) - \mathbf{X}(t))$$

Behauptung

Bei $\gamma(t) := 1/t$ ist immer $\mathbf{w}(t) = \langle \mathbf{X} \rangle_x$

Beweis durch vollständige Induktion $w(0) = 0$

Kap.2.3.2

✓ $\mathbf{w}(t=1) = 0 - \gamma(1)(0 - X) = X = \langle \mathbf{X} \rangle_x$ *Induktionsverankerung*

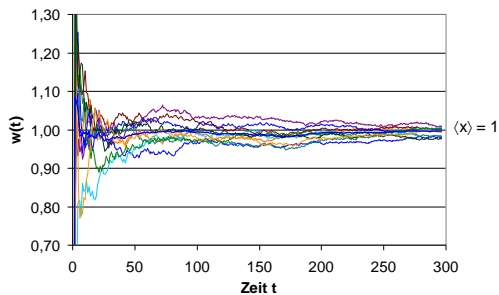
Mit $\mathbf{w}(t-1) = \langle \mathbf{X} \rangle_{t-1} = \sum_{i=1}^{t-1} \mathbf{X}(i)$ *Induktionsvoraussetzung*

✓ gilt $\mathbf{w}(t) = \dots = \langle \mathbf{X} \rangle_t$ *Induktionsschritt*

q.e.d.

Konvergenzverlauf

Iterative Konvergenz



Erwarteter Konvergenzverlauf

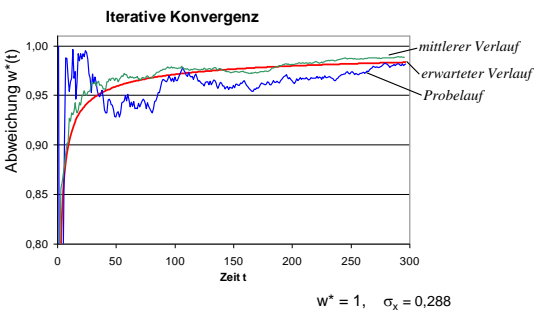
Rechnung

Anhang D.4

- mittl. quadrat. Abweichung
- Erwartungswert aller Verläufe
- Abweichung durch Standardabweichung beschreibbar

$$\langle |w^* - w(t)| \rangle = \sigma_t = \sigma_x / \sqrt{t}$$

Konvergenzverlauf



Stochastisches Lernen

Beispiel Klassentrennung

$$w_i(t) = w_i(t-1) - \gamma(t)(w_i(t-1) - x(t))$$

Behauptung

Bei $\gamma(t) := 1/t$ ist immer $w(t) = \langle x \rangle_x$ *Klassenprototyp*

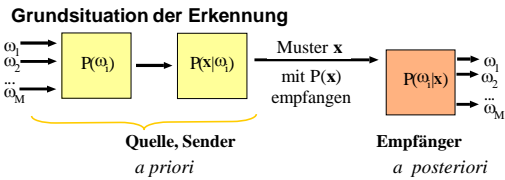
Beweis durch vollständige Induktion $w(0) = 0$

Problem: $\langle x \rangle_x$ ist abhängig von der Klassenentscheidung für x

Assoziativspeicher
Lineare Klassifikation
Lernen linearer Klassifikation
Lernen und Zielfunktion
Stochast. Klassifikation
Lernen in Multilayer-Netzen
Backpropagation-Lernen

Rüdiger Brause: Adaptive Systeme AS-2 WS 2013

Stochastische Musterklassifikation



Notation: $\Omega = \{x\}$, unterteilt in Klassen Ω_k
 $\omega_k = \text{"Klasse } k \text{ liegt vor"}$

Klassifikation
 $\Omega_k: P(\omega_k|x) = \max_j P(\omega_j|x)$ *Bayes-Klassifikation*

Wie erhalten wir $P(\omega_j|x)$?

Rüdiger Brause: Adaptive Systeme AS-2 WS 2013

- 47 -

Stochastische Klassifikation

$P(\omega_j|x) = ?$

Bekannte Quellen

Sei *a-priori* $P(\omega_j)$, $P(x|\omega_j)$, $P(x)$ bekannt
 und $P(x, \omega_j) = P(\omega_j|x)P(x) = P(x|\omega_j)P(\omega_j)$

so ist $P(\omega_j|x) = P(x|\omega_j)P(\omega_j) / P(x)$ mit $P(x) = \sum_j P(x|\omega_j)P(\omega_j)$

Aufgabe !

Unbekannte Quellen

A-posteriori $P(\omega_j|x)$ lernen !

Zielfunktion: Messung der Klassifikationsleistung ?

Rüdiger Brause: Adaptive Systeme AS-2 WS 2013

- 48 -

Klassifikationsleistung

Diagnose-Situation

(Diagnose, Realität)	Name	Wahrscheinlichkeit
$(D(x) = \omega \mid \omega)$	Sensitivität TP	$P_K = P(D(x) = \omega \mid \omega)$
$(D(x) = \neg\omega \mid \omega)$	Ignoranz FP	$P_I = P(D(x) = \neg\omega \mid \omega)$
$(D(x) = \omega \mid \neg\omega)$	Fehlalarm FN	$P_A = P(D(x) = \omega \mid \neg\omega)$
$(D(x) = \neg\omega \mid \neg\omega)$	Spezifität FP	$P_L = P(D(x) = \neg\omega \mid \neg\omega)$

$$P_K + P_I = 1 \quad \text{FP} = \text{FRR} \text{ false rejection rate}$$

$$P_A + P_L = 1 \quad \text{FN} = \text{FAR} \text{ false acceptance rate}$$

Klassifikationsleistung

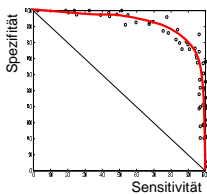
Diagnose-Situation („confusion matrix“)

	ω	$\neg\omega$
$D(x) = \omega$	Sensitivität $P(D(x) = \omega \mid \omega)$ <i>true positive</i>	Fehlalarm $P(D(x) = \omega \mid \neg\omega)$ <i>false negative</i>
$D(x) = \neg\omega$	Ignoranz $P(D(x) = \neg\omega \mid \omega)$ <i>false positive</i>	Spezifität $P(D(x) = \neg\omega \mid \neg\omega)$ <i>true negative</i>

ROC -Kurven von Diagnosesystemen

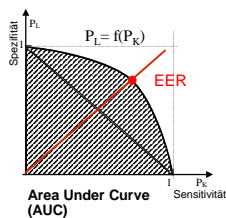
Wechselseit. Abhängigkeit Sensitivität / Spezifität

Beispiel med. Diagnose



Receiver Operating Characteristic (ROC)

Leistung eines Diagnosesystems



Area Under Curve (AUC)

ROC -Kurven von Diagnosesystemen

Aufgabe:

Ex. ein Diagnosesystem mit

$D(x) > c$ Klasse A liegt vor

$D(x) < c$ Klasse A liegt *nicht* vor

Frage:

Wie wird ROC und AUC davon gemessen?

Antwort:

- Für festes c über alle x die Leistung (P_k, P_L) messen, einen Punkt der Grafik einzeichnen
- c variieren, und jeweils Punkt zeichnen
- ROC in die Punkte einpassen, AUC davon berechnen

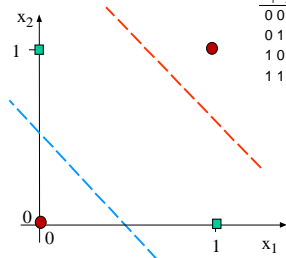
- Assoziativspeicher
- Lineare Klassifikation
- Lernen linearer Klassifikation
- Lernen und Zielfunktion
- Stochast. Klassifikation
- Lernen in Multilayer-Netzen**
- Backpropagation-Lernen

Rüdiger Brause: Adaptive Systeme AS-2 WS 2013

Das XOR-Problem

Aufgabe

Trennung zweier Klassen durch eine Gerade – wie ?



x_1	x_2	$x_1 \oplus x_2$
0	0	0
0	1	1
1	0	1
1	1	0

$$\Omega_0 = \{ \bullet \}$$
$$= \{(0,0), (1,1)\}$$

$$\Omega_1 = \{ \blacksquare \}$$
$$= \{(1,0), (0,1)\}$$

Klassen **nicht**
linear separierbar!

Das XOR-Problem

Lösung

Trennung durch **zwei** Schichten

$$y = -(x_1 \oplus x_2) \quad \text{negiertes XOR}$$

$$= (x_1 \text{ OR } \neg x_2) \text{ AND } (\neg x_1 \text{ OR } x_2)$$

$$y_1 := x_1 \text{ OR } \bar{x}_2$$

$$y_2 := \bar{x}_1 \text{ OR } x_2$$

$$y_{\text{XOR}} := y_1 \text{ AND } y_2$$

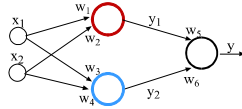
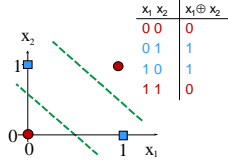
z.B. formale binäre Neuronen

$$S(z \geq s) = 1, S(z < s) = 0$$

$$\Rightarrow w_1 = w_4 = w_5 = w_6 = 1/3$$

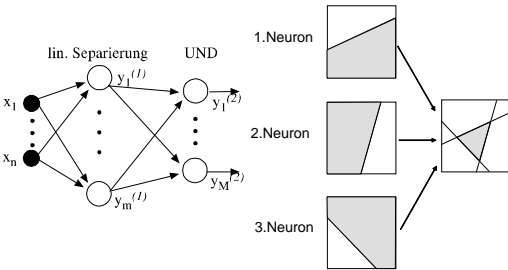
$$w_2 = w_3 = -1/3$$

$$s_1 = s_2 = 0, s = 1/2$$



Multilayer-Klassifikation

Separierung von Klassen



Fähigkeiten der Multilayer-Netzwerke

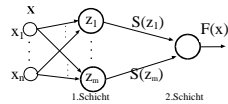
Approximationsnetze

Interpolation anhand von Beispielen (Stützstellen)

Typ. Netz Linearkombinationen von Basisfunktionen $S(\cdot)$

$$F(\mathbf{x}) = \sum_{j=1}^m w_j S(z_j(\mathbf{x}))$$

Sigma-Funktion $F: \mathbb{R}^n \rightarrow \mathbb{R}$



wobei

$$\{ z \mid z(\mathbf{x}) = \mathbf{w}^{(1)T} \mathbf{x} + b \} \quad \text{affine Funktionen } \mathbb{R}^n \rightarrow \mathbb{R}$$

$$\lim_{x \rightarrow \infty} S(x) = 1, \quad \lim_{x \rightarrow -\infty} S(x) = 0 \quad S \text{ ist Quetschfunktion}$$

Fähigkeiten der Multilayer-Netzwerke

Satz Hornik, Stinchcombe, White 1989

Für die Funktionswerte *jeder beliebigen Funktion* $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}^1$ von N Mustern $x^1 \dots x^N$ ex. eine Sigma-Funktion F , so dass für alle Muster x^i mit $i = 1..N$ gilt

$$F(x^i) = f(x^i) \quad \text{Gilt auch für Schicht } \{F_i\}$$

Assoziativspeicher

Satz

Jede beliebige, stetige Funktion $f(x)$ in einem kompakten Intervall ("kompakte Teilmenge des \mathbb{R}^n ") kann **beliebig dicht** (*uniform dicht* im Sinne der L_2 -Norm in der Menge C^n aller stetigen Funktionen und ρ_p -*dicht* in der Menge der Borel-messbaren Funktionen) durch eine Sigma-Funktion $F(x)$ approximiert werden

Anmerkung: Gilt auch für $S =$ stetig, begrenzt, nicht-konstant (RBF)

Fähigkeiten der Multilayer-Netzwerke

Frage:

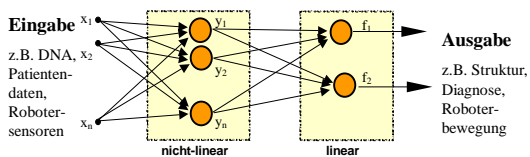
Wieviel Schichten muss ein Netzwerk mindestens haben, um eine beliebige Funktion beliebig gut zu approximieren?

? Antworten:

- eine
- zwei
- drei
- unendlich viele

Mehrschichten-Netze

Fähigkeiten von Mehrschicht-Netzen



- Ein 2-Schichtennetzwerk mit nicht-linearer Ausgabefunktion $S(z)$ kann JEDE beliebige Funktion so genau wie gewünscht approximieren, wenn genügend Neuronen ex.

Neuronenzahl gegeben. Lernalgorithmus=?

Assoziativspeicher
Lineare Klassifikation
Lernen linearer Klassifikation
Lernen und Zielfunktion
Stochast. Klassifikation
Lernen in Multilayer-Netzen

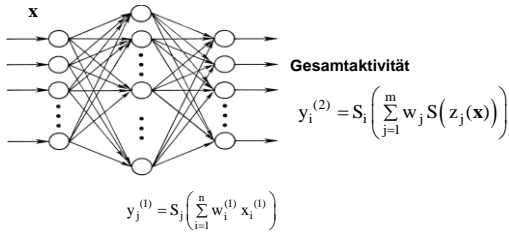
Backpropagation-Lernen

Rüdiger Brause: Adaptive Systeme AS-2 WS 2013

Backpropagation

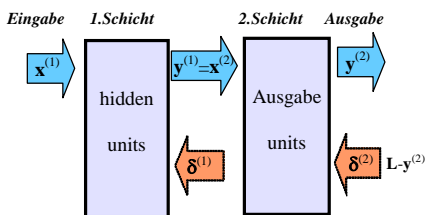
Netzarchitektur und Aktivität

Eingabe hidden units Ausgabe $y_k^{(2)} = S_k \left(\sum_{j=1}^m w_j^{(2)} y_j^{(1)} \right)$



Backpropagation-Grundidee

Netzarchitektur und Lernen



Schichtweise Verbesserung durch Rückführung des Fehlers

Backpropagation-Lernregel letzte Schicht

Lernziel: $R(\mathbf{w}^*) = \min E(y(\mathbf{w}) - L(\mathbf{x}))^2$ min.mittl. quadr. Fehler

$$\mathbf{w}_i(t+1) = \mathbf{w}_i(t) - \gamma \frac{\partial R}{\partial \mathbf{w}_i} \quad \text{Gradienten-Lernregel}$$

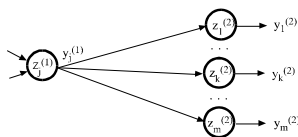
$$\mathbf{w}_{ij}(t+1) = \mathbf{w}_{ij}(t) - \gamma (y_i(\mathbf{w}_{ij}) - L(\mathbf{x})) \frac{\partial S(z_i)}{\partial \mathbf{w}_{ij}} \quad \text{stoch. Approximation}$$

$$\text{mit } \frac{\partial S(z_i)}{\partial \mathbf{w}_{ij}} = S'(z_i) \frac{\partial z_i}{\partial \mathbf{w}_{ij}} = S'(z_i) \frac{\partial}{\partial \mathbf{w}_{ij}} \sum_k \mathbf{w}_{ik} x_k = S'(z_i) x_j$$

$$\text{Mit } \delta_i := - (y_i(\mathbf{w}_{ij}) - L(\mathbf{x})) S'(z_i)$$

$$\text{ist } \Delta \mathbf{w}_{ij}(x) = \gamma \delta_i x_j \quad \text{Delta-Regel}$$

Fehler-Backpropagation



Beeinflussung voriger Schichten $z_1^{(1)} \rightarrow R$

Delta-Regel für Schicht 1, *unabh. von R*

$$\delta_1^{(1)} = - \frac{\partial R_x}{\partial y_1^{(1)}} \frac{\partial y_1^{(1)}}{\partial z_1^{(1)}} = \left(\sum_{k=1}^m \delta_k^{(2)} \mathbf{w}_{k1}^{(2)} \right) S'(z_1^{(1)})$$

Anwendung BP

Gegeben DECTalk

Ausgabe Text \rightarrow Sprache der Fa. Digital Eq. (DEC)
Aufwand 20 PJ für 95% Genauigkeit

Beispiel NetTalk

Sejnowsky, Rosenberg 1986

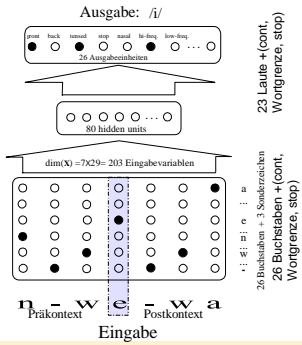
16 CPU-Stunden BP-Training für 98% Genauigkeit

Adaptives Programm statt neu programmieren!

NetTalk: Kodierung

Ausgabekodierung

Binäre Kodierung
der 26 Laute



Eingabekodierung

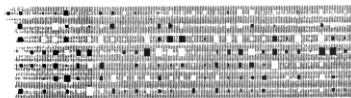
Binäre Kodierung der
29 Buchstaben

Laufenfenster der
Trainingsbuchstaben

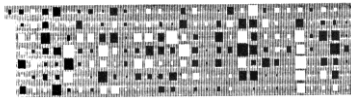
Analyse der Neuronengewichte

Visualisierung der Gewichte

Hinton Diagramm



Gewichte von
Neuron 1



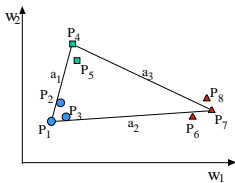
Gewichte von
Neuron 2

□ pos. Gewichte
■ neg. Gewichte

Sinn = ?

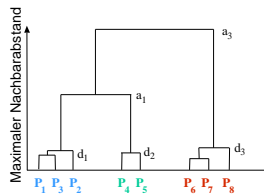
Analyse der Neuronengewichte

Clusteranalyse



$d(x, \text{Nachbar}) < d_N$
 \Rightarrow gleicher Cluster

Hierarchische Clusteranalyse



Dendrogramm

- Sukzessives Zusammenfassung
- Reihenfolge durch Cluster-Abstandsmaß

Analyse der Neuronengewichte

Clusteranalyse

- Clusterung der Muster im Eingaberaum („in vitro“)
- Clusterung der Ausgabewerte bei äquidistanten Testmustern („in vivo“)

→ Funktionsgruppen von Neuronen

Analyse der Neuronengewichte

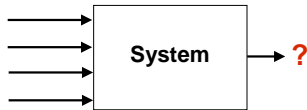
Sensitivitätsanalyse

Aufstellen der Abhängigkeit des Fehlers des Netzes von der Eingabe bzw. den Gewichten.

→ Wichtigkeitsliste der Eingabevariablen

Aber:

- Fehler hängt ab von Signalgrösse
- Normierung d. Signale
- Grosse Gewichte auch bei *random*-Eingabe
- Abhängigkeit von Eingabevar. nicht erfasst



Verbesserungen des BP-Algorithmus

Problem

Das System kann in einem **lokalen** Optimum "stecken" bleiben

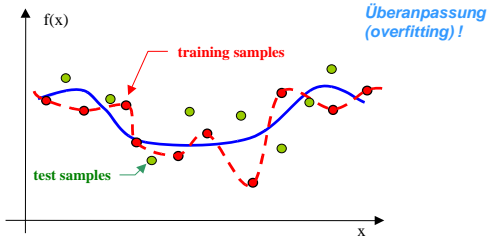
Lösung

- ◆ Gewichte der *hidden units* als Eigenvektoren initialisieren
- ◆ Mehrere Durchgänge mit zufallsveränderten Gewichten
- ◆ Regelmässige Störung der Gewichte & Neulernen
- ◆ Mit kleiner Wahrscheinlichkeit auch gegen das Optimum verändern
- ◆ Sequentieller Netzaufbau, gesteuert durch Kriterium (Ausgabeentropie, Fehler, ...)

Verbesserungen des BP-Algorithmus

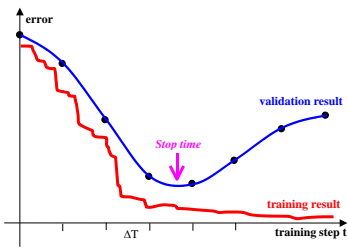
Problem

- Trotz guter Trainingsleistung zeigt der Test **schlechte Ergebnisse**



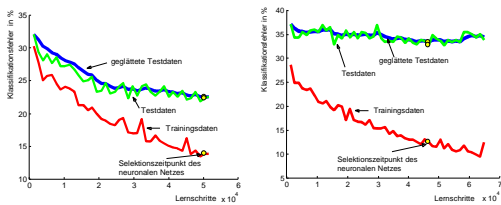
Verbesserungen des BP-Algorithmus

Lösung: Stopped Training



Training und Testen

• Problem: Partition der Daten



zufällige Einteilung der Datenmenge aller Patienten

Aufteilung der Daten nach Patienten

Verbesserungen des BP-Algorithmus

Problem

$$\Delta w_{ij}(x) = \gamma \delta_j x_j = \gamma (\dots) S'(z_j) x_j$$

Die Konvergenz ist relativ **langsam**, besonders bei mehreren Schichten

Problem Ausgabefunktion

Bei Wahl von $S = \text{Fermifunktion}$

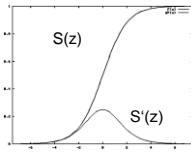
ist die Ableitung eine Glocken-

Funktion mit $S'(-\infty) = 0 = S'(\infty)$

und damit bei sehr großem oder kleinem x

$$\delta(x) = 0$$

⇒ Kein Lernen mehr möglich!



$$S'(z) = \frac{\partial}{\partial z} \frac{1}{1+e^{-z}} = \frac{+1-1+e^{-z}}{(1+e^{-z})^2} = (1-S(z))S(z)$$

Verbesserungen des BP-Algorithmus

Abhilfen für die Ausgabefunktion:

- Andere Lernfunktion wählen, z.B. Ergänzung durch Backpropagation **Trägheitsmoment**

$$\delta(t+1) = \alpha \delta(t) + (1-\alpha) \delta(t-1) \quad \text{z.B. } \alpha = 0.9$$

• Quickprop:

- Addiere 0,1 zu $S'(z) = (1-S)S$
- Ist die Veränderung zu klein, wird sie Null gesetzt
- Ergänzung durch einen **Abklingterm** und **Schrittinterpolation**

$$\Delta w_{ij}(t) = -\gamma(t) R_x(t) \partial w_{ij} - \gamma(t) \alpha w_{ij}(t-1) + \beta(t) \Delta w_{ij}(t-1)$$

- Andere Ausgabefunktion wählen (z.B. Sinus)

Verbesserungen des BP-Algorithmus

Problem

Die Konvergenz ist relativ **langsam**, besonders bei mehreren Schichten

Lösung

- Abhilfe bei Ausgabefunktion
- Wahl von $\gamma(t)$ als Hauptdiagonal-Matrix
- Änderung von $\gamma(t)$ bei Vorzeichen des Gradienten
- Lernen der Einzelschichten, falls möglich (z.B. zuerst EV bei hidden units verwenden und 2. Schicht lernen, dann 1. Schicht)
- Andere Zielfunktion wählen (Information statt erwart. quadr. Fehler MSE)

Wechsel der Zielfunktion

Zielfunktion Information

$$P_{kj} = P(\omega_j | x^k)^{L_j} (1 - P(\omega_j | x^k))^{1-L_j}$$

Wahrsch. für Klassifik. von
Muster k in Klasse j bei
Lehrerurteil L

$$P_k = \prod_{j=1}^M P_{kj} = \prod_{j=1}^M P(\omega_j | x^k)^{L_j} (1 - P(\omega_j | x^k))^{1-L_j}$$

Wahrsch. bei M Entscheidng., Muster k richtig zu klassifizieren

DEF $R_x := I(x^k) = -\log P_k$ *Zielfunktion* *Kap. 2.6.9*

$$= -\sum_{j=1}^M L_j \log P(\omega_j | x^k) + (1-L_j) \log (1 - P(\omega_j | x^k)) \quad y_j = P(\omega_j | x^k)$$

$$= -\sum_{j=1}^M L_j \log y_j + (1-L_j) \log (1-y_j) \quad \Rightarrow \delta^{(2)} = (y-L)$$

Wechsel der Zielfunktion

Beispiel: Klassifikation von Phonemen

MSE 68% ok

Information 78% ok

