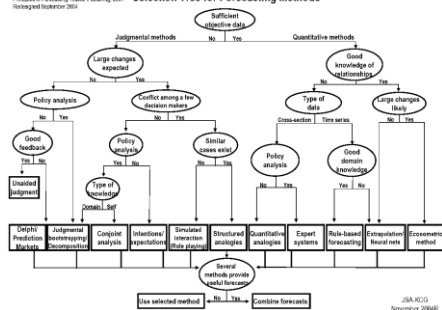


## Kapitel 2: Black-Box-Modellierung

## Rüdiger Brause

### Selection Tree for Forecasting Methods



R.Brause, Adaptive Modellierung: Kap.2 Black-Box-Modellierung

304R

Beispiel		Individuelle Produktion und Importe MW, EUR											
		Land der jeweiligen Herkunft											
		China v. Mexiko	China v. USA	USA v. China	USA v. Mexiko	China v. USA	China v. Mexiko	USA v. China	USA v. Mexiko	China v. USA	China v. Mexiko	USA v. China	USA v. Mexiko
		1	2	3	4	5	6	7	8	9	10	11	12
Aufgaben	1. "Glocke" nach Deutschland (China 100, Mexiko 120)	1,6	0,9	0,5	-	0,9	2,1	34,9	0,9	1,6	0,9		
	2. Energie und Wasser (China 100, Mexiko 120)	1,6	0,7	25,7	7,0	4,0	3,3	2,9	2,6	7,6	2,6		
	3. Wein (China 100, Mexiko 120)	1,6	0,9	0,5	-	0,9	2,1	34,9	0,9	1,6	0,9		
	4. Weizen (China 100, Mexiko 120)	0,6	1,0	4,7	4,0	3,0	3,0	1,2	12,1	3,3	0,5		
	5. Autos (China 100, Mexiko 120)	0,6	1,0	4,7	4,0	3,0	3,0	1,2	12,1	3,3	0,5		
	6. Tomaten, Gurken und Kohl (China 100, Mexiko 120)	0,3	0,6	4,7	3,0	3,0	3,0	1,2	12,1	3,3	0,5		
	7. Nahrungsmittel und Elektronik (China 100, Mexiko 120)	4,0	0,9	2,7	0,8	0,1	0,6	49,7	0,6	0,9	12,9		
	8. Haushaltswaren (China 100, Mexiko 120)	0,6	2,1	1,9	1,9	1,9	1,9	0,7	0,6	4,7	0,3	21,9	
	9. Haushaltswaren und Elektronik (China 100, Mexiko 120)	3,9	3,1	15,6	11,4	37,9	12,9	12,6	3,9	13,6	15,6		
	10. Haushaltswaren und elektronische Haushaltswaren (China 100, Mexiko 120)	4,0	0,9	29,7	10,2	31,2	10,6	16,6	33,2	30,9	29,2		
11. Haushaltswaren und elektronische Haushaltswaren (China 100, Mexiko 120)	1,0	0,6	1,6	0,7	0,7	0,6	0,9	0,9	0,6	0,6			
12. Haushaltswaren und elektronische Haushaltswaren (China 100, Mexiko 120)	0,3	4,1	1,9	0,7	1,7	3,6	1,9	1,1	0,1	16,2			
Summe		22,9	26,7	100,3	63,3	206,0	10,9	82,7	124,0	10,9			

R.Brause, Adaptive Modellierung: Kap.2 Black-Box-Modellierung

## Lineare Modellierung

Nichtlin. Modellierung: MLP

Nichtlin. Modellierung: RBF

Eigenschaften von NN

---

---

---

---

---

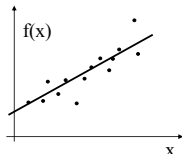
---

---

---

## Lineare Approximation

•  $m$  Messungen



• Modellierung als Gerade

$$y = f(x) = y_0 + ax$$

• Beispiel: Ökonomie

Konsum  $y = f(\text{Einkommen } x)$   
= Konsumsokkel +  $a \cdot \text{Einkommen}$

---

---

---

---

---

---

---

---

## Lineare Approximation - rauschfrei

• Parameterbestimmung

2 Messwerte  $y_1(x)$ ,  $y_2(x)$  reichen aus für  $a$ ,  $y_0$   
RECHNUNG

•  $x, y, y_0$  sind  $n$ -dim Vektoren,  $a$  eine  $n \times n$ -Matrix  
 $m = n^2 + n$  Parameter  $\Rightarrow m$  Datenpunkte bei  $n$  Variablen  
nötig

---

---

---

---

---

---

---

---

## Lineare Approximation - verrauscht

- Modellierung als verrauschte Gerade

$$y_t = y_0 + ax_t + u_t$$

- Parameterbestimmung **Rechnung**  
aus Varianz und Kovarianz

- Parameterbestimmung **Rechnung**  
mittels Gauß-Methode

$$R(a) = \langle (y_t - f(a))^2 \rangle_t = \min_a$$

## Lineare Approximation - verrauscht

- Mehrere Variable *Multiple Regression*

$$y(t) = f(x_1, x_2, \dots, x_k) = a_0 + a_1 x_1 + a_2 x_2 + \dots + a_k x_k + u$$

Zusammenfassung zu

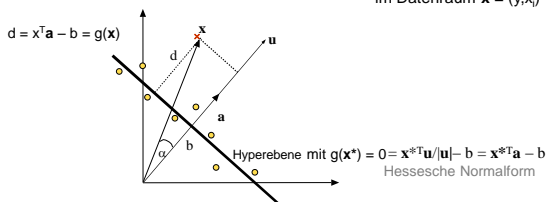
$$y_t = (1, x_1(t), \dots, x_k(t)) (a_0, a_1, \dots, a_k)^T + u_t = \mathbf{x} \mathbf{a}^T + u_t$$

$$\mathbf{y} = (y_1, \dots, y_T)^T = \mathbf{X} \mathbf{a}^T + \mathbf{u} \quad T \text{ Zeilen für } T \text{ Messungen, } \text{rang}(\mathbf{X}) = k+1$$

- Parameterbestimmung **Rechnung**

## Lineare Approximation - verrauscht

- Minimaler mittlerer Abstand zur Gerade (Hyperebene)  
im Datenraum  $\mathbf{x} = (y, x_i)$



$$\text{TLMSE} = R(\mathbf{a}, \mathbf{b}) = \langle d^2 \rangle$$

**Rechnung:** Minimum des TLMSE (Kap.2.2)

## TLMSE - Parameterschätzung

- **Vorhanden:** Messungen  $\mathbf{x} = (y, \mathbf{x}^T)$
- **Gesucht:** Eigenvektor von  $\mathbf{C}_{\mathbf{x}\mathbf{x}}$  mit min.  $\lambda$

- **1. Lösung:** Fixpunktalgorithmus für EV

$\mathbf{a}(t+1) = \mathbf{C}\mathbf{a}(t)$ ,  $|\mathbf{a}| = 1 \rightarrow$  EV mit max. EW  
Neuer Eingaberaum  $\mathbf{x}' = \mathbf{x} - \mathbf{a}\mathbf{a}^T\mathbf{x}$ ,  $\mathbf{C}_{\mathbf{x}'\mathbf{x}'}$  bilden, nächsten EV lernen.

- **2. Lösung:** Anti-Hebb-Lernalgorithmus

$\mathbf{a}(t) = \mathbf{a}(t-1) - \mathbf{x}(\mathbf{a}^T\mathbf{x})$ ,  $|\mathbf{a}| = 1$  *Anti-Hebb Lernregel*

## Lineare Approximation - Grenzen

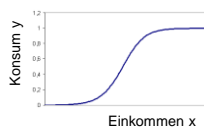
- **Grenzen des linearen Modells**

- linearer Zusammenhang der Daten ? Linearisieren!
- Parameter sind konstant (oder zeitabhängig?)
- Einflüsse von  $x_1, x_2, \dots$  sind nur eingebildet ?

## Lineare Approximation - Grenzen

- **Nichtlineare Modellierung: Beispiel Ökonomie**

Grossverdiener: Wer alles hat, kauft nichts mehr.  
Sigmoidaler Zusammenhang ist wahrscheinlicher:



## Linearisierung

### Linearisieren von nichtlin. Messungen

- $y = y_0 + ax^{1/2} + u$   
 $z = x^{1/2} \rightarrow y(z) = y_0 + az + u$
- $y = ax_1^\alpha \cdot x_2^\beta \cdot e^u$  („Cobb-Douglas-Produktion“)  
 $z = \ln(y) \rightarrow z(x) = \ln a + \alpha \ln x_1 + \beta \ln x_2 + u$   
 $= a_0 + a_1 x_1' + a_2 x_2' + u$
- nicht-lin. Kontext, z.B. „männlich“ bei  $y = y_0 + ax_1 + u$   
 $x_2 = 1$  bei „männlich“, sonst null  
 $\rightarrow y = y_0 + ax_1 + bx_2 + \dots + u$

R.Brause, Adaptive Modellierung: Kap.2 Black-Box-Modellierung

- 213 -

---

---

---

---

---

---

---

---

---

---

## Lineare Approximation

### Parameter konstant ?

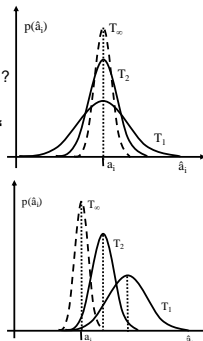
- Schätzung  $\hat{a}_i$  liegt vor.
- Parameter  $a_i \neq 0$  ? Oder nur Störung ?

### Kriterium „erwartungswertreu“ (unbiased)

$$\langle \hat{a}_i(t) \rangle_t = a_i$$

### Kriterium „konsistent“

$$\lim_{t \rightarrow \infty} \text{Prob} \left( |\hat{a}_i(t) - a_i| < \varepsilon \right) = 1$$



R.Brause, Adaptive Modellierung: Kap.2 Black-Box-Modellierung

- 214 -

---

---

---

---

---

---

---

---

---

---

## Lineare Approximation - Störterm

### Störterm $\langle u_t \rangle = u_0 \neq 0$

Was tun ? Mittelwert in Konstante verschieben:

$$y_t = (a_0 + u_0) + a_1 x_1(t) + a_2 x_2(t) + \dots + a_k x_k(t) + (u_t - u_0)$$

$$= \hat{a}_0 + a_1 x_1(t) + a_2 x_2(t) + \dots + a_k x_k(t) + \hat{u}_t$$

### Forderung: Keine Korrelation des Störterms mit den Variablen x !

$$\text{cov}(u_t, u_{t'} | \mathbf{x}_t) = 0 \quad \forall t, t' = 1, \dots, T \quad t \neq t'$$

R.Brause, Adaptive Modellierung: Kap.2 Black-Box-Modellierung

- 215 -

---

---

---

---

---

---

---

---

---

---

## Lineare Approximation - Multikollinearität

### Test der Abhängigkeit der Variablen untereinander

- Bilde Korrelationskoeffizienten

$$r_{ij} = \frac{\text{cov}(x_i, x_j)}{\sqrt{\text{var}(x_i)} \sqrt{\text{var}(x_j)}} = \frac{\text{cov}(x_i, x_j)}{\sigma_i \sigma_j}$$

- Bilde Bestimmtheitsmaß

$$R^2 = \frac{\langle (\hat{y}_i - \bar{y})^2 \rangle}{\langle (y_i - \bar{y})^2 \rangle} \leq \frac{\langle \hat{u}_i^2 \rangle}{\langle (y_i - \bar{y})^2 \rangle} + \frac{\langle (\hat{y}_i - \bar{y})^2 \rangle}{\langle (y_i - \bar{y})^2 \rangle} = 1$$

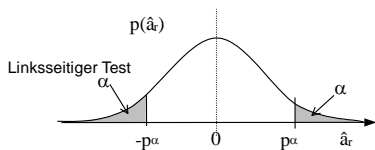
- Teste: Gilt  $r_{ij}^2 \geq R^2$  ?

JA: **Multikollinearität liegt vor!**

## Lineare Approximation

### Test der Einflüsse: Ist Variable $x_r$ nötig? Ist der Parameter $a_r = 0$ ?

Nullhypothese  $H_0 : a_r = 0$



Rechtsseitiger Test:

Ist  $\text{Prob}(\hat{a}_r > p_\alpha) \leq \alpha$  ?

JA:  $H_0$  ok.

Verteilung der beobachteten Parameterwerte

Hypothesentest auch für  $a_r = s$  brauchbar!

## Black-Box- Modellierung

### Lineare Modellierung

### Nichtlin. Modellierung: MLP

### Nichtlin. Modellierung: RBF

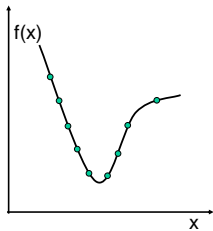
### Eigenschaften von NN

## NichtLineare Approximation

### Polynomansatz ( Taylorentwicklung)

$$f(x) = f(x_0) + (x-x_0)f^{(1)}(x_0) + \frac{(x-x_0)^2}{2} f^{(2)}(x_0) + \frac{(x-x_0)^3}{6} f^{(3)}(x_0) + \dots$$

$$= A + B(x-x_0) + C(x-x_0)^2 + D(x-x_0)^3 + \dots$$



Direkte Bestimmung  
der Parameter aus den  
Messwerten

---

---

---

---

---

---

---

---

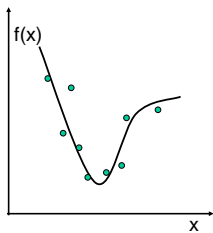
---

---

## NichtLineare Approximation

### Polynomansatz mit Messfehlern

$$f(x) = A + B(x-x_0) + C(x-x_0)^2 + D(x-x_0)^3 + \dots$$



Indirekte Bestimmung der  
Parameter aus den  
Messwerten:  
**Regression n-ter Ordnung,**  
z.B. mit kleinstem quadr. Fehler

---

---

---

---

---

---

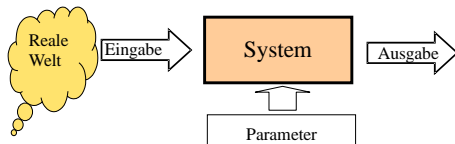
---

---

---

---

## Nichtlineare Approximation



### Einsatzarten

- Adaptive Schätzung von Prozeßparametern  
Nicht-lin. Reaktionen, Produktionsoptimierung,...
- Adaptive Kontrolle und Regelung  
Landekontrollsysteme, Roboterkontrolle,...
- Adaptive Klassifikation  
Qualitätskontrolle, med. Diagnose, Bonitätsprüfung,...

---

---

---

---

---

---

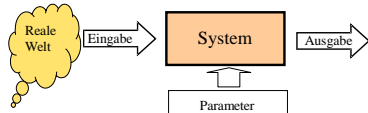
---

---

---

---

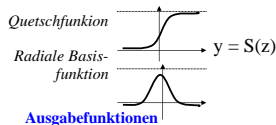
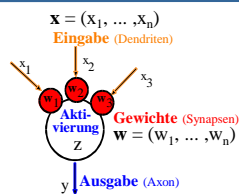
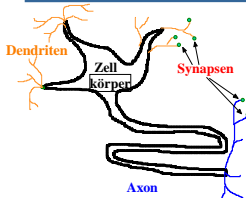
## Nichtlineare Approximation



### Einsatzgebiete

- Echtzeitreaktionen  
z.B. Stahlwalzstraßen, Flugzeugsteuerung,...
- Analytisch unbekannte Abhängigkeiten  
z.B. Polymerchemie, DNA-Schätzungen, ...
- Analytisch nicht zugängige Abhängigkeiten  
z.B. psychische Faktoren, ergebnisverändernde Messungen,...
- Analytisch nur unter großem Aufwand bearbeitbare, hochdimensionale Gesetzmäßigkeiten  
z.B. Wechselkursabhängigkeiten,...
- Statistische Analysen durch untrainierte Benutzer (!)

## Was sind Neuronale Netze ?



$$\mathbf{x} = (x_1, \dots, x_n)$$

Eingabe (Dendriten)

$$\mathbf{w} = (w_1, \dots, w_n)$$

Gewichte (Synapsen)

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

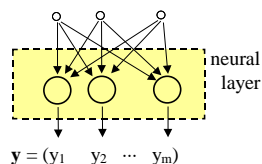
Aktivierung

$$y = S(z)$$

## Was sind Neuronale Netze ?

### DEF Schicht

$$\mathbf{x} = (x_1 \quad x_2 \quad \dots \quad x_n)$$



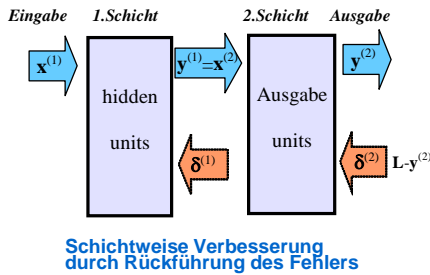
### Lineare Schicht

$$\mathbf{y} = \left( \sum_i w_{1i} x_i, \dots, \sum_i w_{mi} x_i \right)^T = \mathbf{W} \cdot \mathbf{x} \quad \text{Matrixmultiplikation}$$



## Backpropagation-Grundidee

### Netzarchitektur und Lernen



R. Brause, Adaptive Modellierung: Kap. 2 Black-Box-Modellierung

- 25 -

## Anwendung - NetTalk

Sejnowsky-Rosenberg 1986

- Automatisches System Text -> Sprache
- Vorläufer: DECTalk, 20 Mann-Jahre, 95% Genauigkeit
- NetTalk: 15 CPU-Stunden, 98% Genauigkeit
- Eingabe: Texte

R. Brause, Adaptive Modellierung: Kap. 2 Black-Box-Modellierung

- 26 -

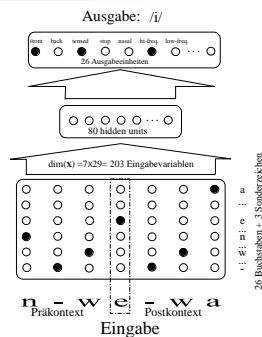
## Anwendung - NetTalk

### Architektur

**Eingabe:** 29 Zeichen, binär  
(26 Buchstaben+ 3 Sonderzeichen:  
Pkt., Wortgrenze),  
7 Buchstaben als Kontext.

**Hidden units:** 80 Stück

**Ausgabe:** 26 Merkmale, binär  
(23 Laut- und drei Artikulations-  
merkmale: Continuation,  
Wortgrenze, Stop)



R. Brause, Adaptive Modellierung: Kap. 2 Black-Box-Modellierung

- 27 -

## Anwendung - NetTalk

### Training

- protokollierte, von Kindern gesprochenen Sätze,
- zufallsmäßig eingegebene Worte eines Wörterbuchs aus 20.000 Einträgen

Eingabe x: Buchstabe eines Worts im Kontext  
Lehrervorgabe L(x): Phonologische Transkription  
Einfügen eines Sonderzeichens „Continuation“, wenn Buchstabe nicht gesprochen wird (!)

## Anwendung - NetTalk

### Ergebnisse: 3 Phasen des Sprachlernens

- Zuerst wurden die Konsonanten und Vokale als Klassen getrennt. Innerhalb der Klassen blieben die Phoneme aber noch gemischt, so dass sich die Laute wie "Babbeln" anhörten.
- Dann wurden die Wortgrenzen als Merkmale entwickelt, so dass "Pseudoworte" erkennbar wurden.
- Zuletzt, nach ca. 10 Durchgängen pro Wort, entstand eine verständliche Sprache, die sich mit fortlaufender Erfahrung weiter verbesserte

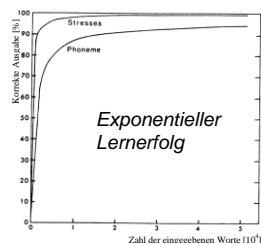
## NetTalk: Training

### Training

- transkribiertes Wörterbuch 20.000 Einträge
- Protokollierte Kindersätze

### Ergebnis

- Trennung der Konsonanten von Vokalen („Babbeln“)
- Entwicklung der Wortgrenzen („Pseudoworte“)
- Verständliche Sprache (10x Training pro Wort)



## Backpropagation-Grundidee

### • Zielfunktion minimieren

Ziel = minimaler quadrat. Fehler

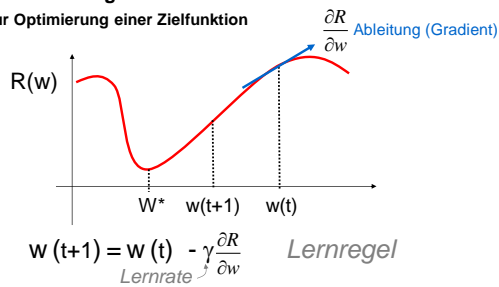
$$R(\mathbf{w}) = \langle (L(x) - y^{(2)}(x))^2 \rangle_x = \min_{\mathbf{w}}$$

Wie erreichen?

## Lernen mit Zielfunktionen

### Gradientenalgorithmus

zur Optimierung einer Zielfunktion



## Stochastisches Lernen

Lernen mit Zielfunktion  $R(\mathbf{w}) = \langle R_x(\mathbf{w}, x) \rangle_x$

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

wird ersetzt durch

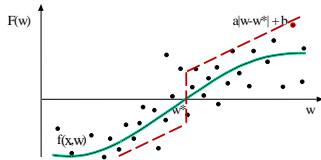
Lernen mit **stochast. Zielfunktion**  $R_x(\mathbf{w}, x)$

$$\mathbf{w}(t) = \mathbf{w}(t-1) - \gamma(t) \nabla_{\mathbf{w}} R_x(\mathbf{w}(t-1), \mathbf{x}(t)) \quad \text{stochastisches Lernen}$$

Wieso darf man das?

## Stochastische Approximation

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



**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*

$$\mathbf{w}(t) = \mathbf{w}(t-1) - \gamma(t) f(\mathbf{w}(t-1), \mathbf{x}(t))$$

*stochastische  
Approximation*

## Stochastische Approximation

JOHANN WOLFGANG GOETHE  
UNIVERSITÄT  
FRANKFURT AM MAIN

- ## 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.  $\sigma^2(F(w)) = \langle (F(w) - f(x, w))^2 \rangle_x < \infty$
  - $\gamma(t)$  **verschwindet**,  $\gamma(t) \rightarrow 0$
  - $\gamma(t)$  wird **nicht zu schnell klein**,  $\sum_{i=1}^{\infty} \gamma(t_i) = \infty$
  - $\gamma(t)$  wird **nicht zu groß**,  $\sum_{i=1}^{\infty} \gamma(t_i)^2 < \infty$

**Dann ex.**

$\lim_{t \rightarrow \infty} \langle (w(t) - w^*)^2 \rangle = 0$	mittl. quadr. Konvergenz <i>Robbins-Monro</i>
$P(\lim_{t \rightarrow \infty} w(t) = w^*) = 1$	<i>Blum</i>

## Backpropagation-Lernregel letzte Schicht

JOHANN WOLFGANG GOETHE  
UNIVERSITÄT  
KUFURT AM MAIN

**Lernziel:**  $R(w^*) = \min \langle (y(x,w) - L(x))^2 \rangle_x$  min.mittl. quadr. Fehler

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

$$\mathbf{w}_i(t+1) = \mathbf{w}_i(t) - \gamma (y(\mathbf{w}_i) - L(\mathbf{x})) \frac{\partial S(z)}{\partial w_i} \quad \text{stoch. Approximation}$$

### Rechnung: Ableitung der Zielfunktion

## Backpropagation-Lernregel letzte Schicht

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

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

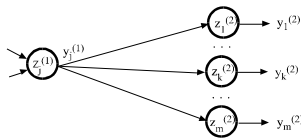
$$w_i(t+1) = w_i(t) - \gamma (y(w_i) - L(x)) \frac{\partial S(z)}{\partial w_i} \quad \text{stoch. Approximation}$$

$$\text{mit } \frac{\partial S(z)}{\partial w_i} = \frac{S'(z) \partial z}{\partial w_i} = \frac{S'(z) \partial}{\partial w_i} \sum_j w_j x_j = S'(z) x_i$$

$$\delta_i := - (y(w_i) - L(x)) S'(z)$$

$$\Delta w_{ij}(x) = \gamma \delta_i x_j \quad \text{Delta-Regel}$$

## Fehler-Backpropagation



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

Delta-Regel für Schicht 1

$$\delta_i^{(1)} = \frac{-\partial R_x}{\partial y_i^{(1)}} \frac{\partial y_i^{(1)}}{\partial z_i^{(1)}} = \left( \sum_{k=1}^m \delta_k^{(2)} w_{ki}^{(2)} \right) S'(z_i^{(1)})$$

## Online vs Offline-Lernen

### ONLINE-Learning:

**WHILE NOT** Abbruchbedingung erfüllt:

Delta := 0

**FORALL** Trainingsmuster x

berechne Delta(W(x))

W(t) := W(t-1) + Delta // Lernen mit jedem Muster

**END FOR**

**END WHILE**

## Online vs Offline-Lernen

### OFFLINE-Learning:

**WHILE NOT** Abbruchbedingung erfüllt:

GesamtDelta := 0

**FORALL** Trainingsmuster x

berechne Delta(W(x))

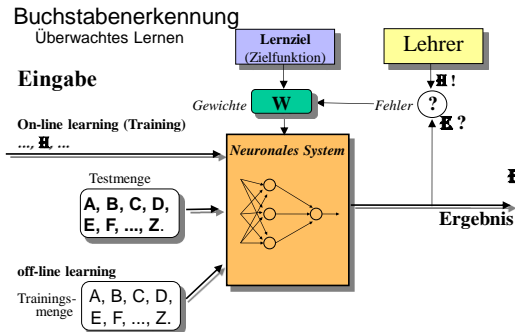
GesamtDelta := GesamtDelta + Delta(W(x))

**END FOR**

W(t) := W(t-1) + GesamtDelta // Lernen am Schluss!

**END WHILE**

## Arten des Lernens: Beispiel



## Frage: Offline- oder Online-Lernen?

### Backpropagation- Code

```
REPEAT (* jeweils einen Trainingszyklus *)
  dw1:=0.0; dw2:=0.0; y:=0.1;
  REPEAT (* Für alle Trainingsmuster im PatternFile *)
    Read( PatternFile,x1,L); (*Einlesen der Eingabe, Ausgabe *)
    (* Ausgabe errechnen *)
    FOR i:=1 TO p DO x2[i]:= S(z(w1[i],x1)) END (* Für hidden units *)
    FOR i:=1 TO m DO (* Für alle Ausgabeneuronen *)
      y2[i]:= S(z(w2[i],x2))
      d2[i]:= -(y2[i]-L[i])*(1-y2[i])*y2[i] (* (y-L)/(1-S)S *)
    END
    FOR i:=1 TO m DO (* Gewichtsveränderungen in 2. Schicht *)
      FOR j:=1 TO p DO
        dw2[i,j] := dw2[i,j] + y*d2[i]*x2[j]
      END;
    END
    FOR i:=1 TO p DO (* Gewichtsveränderungen in 1. Schicht *)
      FOR j:=1 TO n DO (* Für alle Eingabemusterkomp. *)
        dw1[i,j] := dw1[i,j]+y*SumProd(i,m,d2,w2)*(1-x2[i])*x2[i]*x1[j]
      END;
    END
  UNTIL ( EOF( PatternFile) )
  w1:=w1+dw1; w2:=w2+dw2; (* Korrektur der Gewichte *)
UNTIL Fehler_klein_genug
```

## Lernen - Probleme

- Die Konvergenz ist relativ *langsam*, besonders bei mehreren Schichten
- Das System kann in einem *lokalen Optimum* "stecken" bleiben
- Trotz guter Trainingsleistung zeigt der Test *schlechte Ergebnisse*

---

---

---

---

---

---

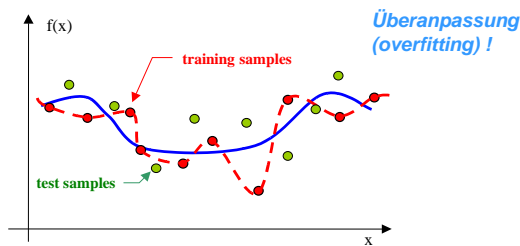
---

---

## Verbesserungen des BP-Algorithmus

### Problem

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




---

---

---

---

---

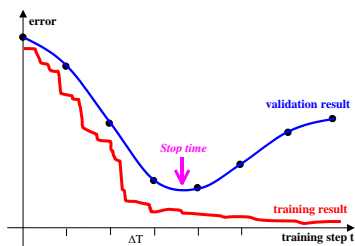
---

---

---

## Verbesserungen des BP-Algorithmus

### Lösung: Stopped Training




---

---

---

---

---

---

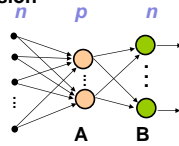
---

---

## Initialisierung der Neuronengewichte

### Beispiel Informationskompression

- **Lin. Approximation**  
(1. Glied Taylorentwicklung)  
Beispiel: n-p-n Netz Kodierer  
 $y = B_{p \times n} A_{n \times p} x$

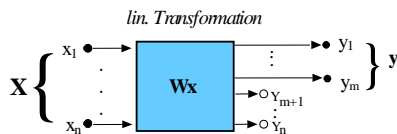


- **Min. quadr. Fehler** bei globalem Minimum

⇒ **A = ?**

## Transformation mit minimalem MSE

### Allgemeine Situation



$$\min_w R(W) = \min \langle (x - \hat{x})^2 \rangle \quad \text{least mean squared error (LMSE)}$$

**Wann minimal ?**

## Transformation mit minimalem MSE

### Minimaler Rekonstruktionsfehler

$$\min_w R(W) = \min \langle (x - \hat{x})^2 \rangle \quad \text{least mean squared error (LMSE)}$$

$$x = \sum_{i=1}^m y_i w_i + \sum_{i=m+1}^n y_i w_i \quad \hat{x} = \sum_{i=1}^m y_i w_i + \sum_{i=m+1}^n c_i w_i \quad y_i = x^T w_i$$

- Was ist die beste Schätzung für die Konstanten  $c_i$ ?  
 $\min R(c_i) = ?$  **Rechnung!**

- Bei welchen Basisvektoren  $w_i$  ist der Fehler minimal ?  
 $\min R(w_i) = ?$  **Rechnung!**



## Analyse der Neuronengewichte

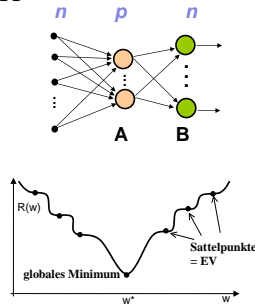
### Hauptkomponentenanalyse

Lin. Approximation  
(1. Glied Taylorentwicklung)

• **A** besteht aus **Eigenvektoren** der  
**Kovarianzmatrix**

$$C_{xx} = \langle (\mathbf{x} - \langle \mathbf{x} \rangle) (\mathbf{x} - \langle \mathbf{x} \rangle)^T \rangle$$

$$(C_{ij}) = \left( \frac{1}{N} \sum_{k=1}^N (x_i^k - \langle x_i \rangle) (x_j^k - \langle x_j \rangle) \right)$$



R.Brause, Adaptive Modellierung: Kap.2 Black-Box-Modellierung

-49-

### Normierung der Eingangsvariablen

• **Problem:** unterschiedliche Skalierungen

- z.B.
- $x_1$  : Weg in [m]
  - $x_2$  : Grösse in [cm]
  - $x_3$  : Gewicht in [kg]
  - $x_4$  : Farbwert in [RGB-Zahl]

Normierung aller numerischen Werte auf gleiche Skala !

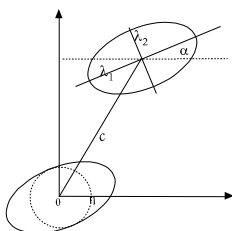
- Gleichen Mittelwert, etwa  $\langle x \rangle = 0$
- Gleiche Varianz  $\sigma = 1$

R.Brause, Adaptive Modellierung: Kap.2 Black-Box-Modellierung

-50-

### Normierung der Eingangsvariablen

• **Abstände bilden nur mit normierten Variablen**



**Mahalanobis-Abstand**

$$d^2 = (\mathbf{x} - \mathbf{c})^T \mathbf{C}^{-1} (\mathbf{x} - \mathbf{c})$$

$$= (\mathbf{x} - \mathbf{c})^T \mathbf{M}^T \mathbf{M} (\mathbf{x} - \mathbf{c})$$

Entspricht einer Skalierung,  
Drehung, Verschiebung mit

$$\mathbf{x} \rightarrow \mathbf{Mx} = \mathbf{SDVx} \quad \text{mit } \mathbf{x} \rightarrow (\mathbf{x}^T, 1)^T$$

$$\mathbf{S} = \begin{pmatrix} s_1 & 0 \\ 0 & s_2 \end{pmatrix}, \quad \mathbf{D} = \begin{pmatrix} \cos \alpha & -\sin \alpha \\ \sin \alpha & \cos \alpha \end{pmatrix}$$

$$\mathbf{V} = \begin{pmatrix} 1 & 0 & -c_1 \\ 0 & 1 & -c_2 \end{pmatrix}$$

R.Brause, Adaptive Modellierung: Kap.2 Black-Box-Modellierung

-51-

## Verbesserungen des BP-Algorithmus

### Problem

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

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

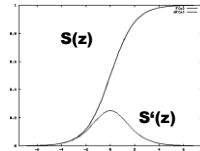
### Problem Ausgabefunktion

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!**



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

## Problem Ausgabefunktion

### Abhilfen:

- Andere Ausgabefunktion wählen (z.B. Sinus)
- Andere Zielfunktion wählen (Information statt quadr.Fehler)
- Andere Lernfunktion wählen, z.B. Ergänzung durch Backpropagation *Trägheitsmoment*

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

## Black-Box- Modellierung

### Lineare Modellierung

### Nichtlin. Modellierung: MLP

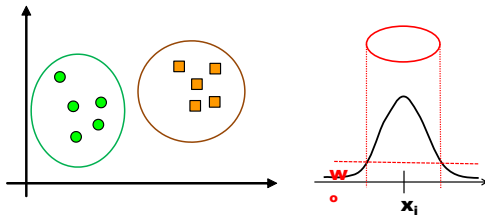
### Nichtlin. Modellierung: RBF

### Eigenschaften von NN

## Klassifikation und RBF

Motivation: lokale Cluster-Klassenbildung

$$\Omega_i = \{ \mathbf{x} \mid S(|\mathbf{x} - \mathbf{x}_i|) > w_0 \}$$



Rüdiger Brause: Adaptive Systeme, Institut für Informatik, WS 2013/14

- 55 -

---

---

---

---

---

---

---

---

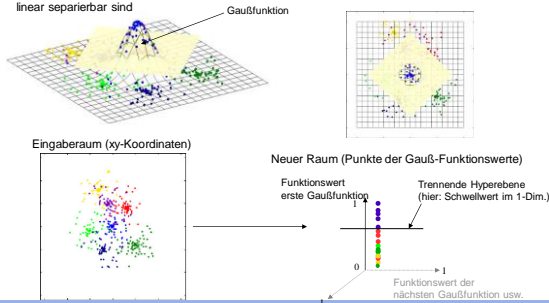
---

---

## Klassifikation und RBF

Bessere Klassifikation durch gekrümmte Entscheidungsgebiete

Idee: Nichtlineare Abbildung in neuen Raum, in dem die Muster (mit höherer Wahrscheinlichkeit) linear separierbar sind



R. Brause, Institut für Informatik

- 56 -

---

---

---

---

---

---

---

---

---

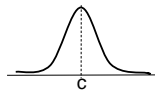
---

## Radiale Basisfunktionen

Definition **Glockenfunktion**

Funktion  $S_G$  mit den Eigenschaften

- $S_G(z) \geq 0$ ,  $S_G(-\infty) = S_G(\infty) = 0$ ,
- $0 < \int_{-\infty}^{\infty} S_G(x) dx < \infty$
- Es ex. ein  $c > 0$  mit  $S_G(z)$  nicht anwachsend  $\forall z \in [c, \infty)$ ,  
nicht abfallend  $\forall z \in (-\infty, c)$



Also ist  $S_G(c)$  globales Maximum.

Rüdiger Brause: Adaptive Systeme, Institut für Informatik

- 57 -

---

---

---

---

---

---

---

---

---

---

## RBF maximaler Information

Welche Basisfunktionen hat maximale Information ?

$$H(p^*) = \max_p H(p(x)) \quad x \in \mathbb{R}, \quad p^*(x) = ?$$

NB1:  $\int p(x) dx = 1$  oder  $g_1(x) := \int p(x) dx - 1 = 0$

NB2:  $\sigma^2 = \langle x^2 \rangle = \int_{-\infty}^{\infty} p(x) x^2 dx$  oder  $g_2(x) := \int_{-\infty}^{\infty} p(x) x^2 dx - \sigma^2 = 0$

Ansatz **Lagrange-Funktion**

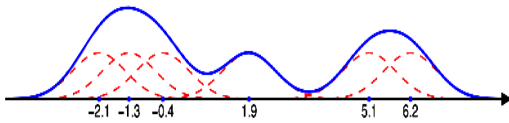
$$L(p, \mu_1, \mu_2) := H(p) + \mu_1 g_1(p) + \mu_2 g_2(p)$$

$$\frac{\partial L(p^*)}{\partial p} = 0, \quad \frac{\partial L(\mu_i)}{\partial \mu_i} = 0 \quad (\text{Rechnung Kap.5.2})$$

Ergebnis  $p^*(x) = A \exp(-x^2/2\sigma^2)$  **Gauß'sche Glockenkurve**

## Parzen Window - Methode

Approximation durch Überlagerung von Basisfunktionen



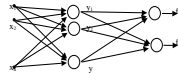
$$\rho(x) = \frac{1}{N} \sum_{i=1}^N W(x - x_i) \quad W(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-x^2/2\sigma^2}$$

Perfekte Approximation bei abnehmender Breite  $\sigma$ , wobei

$$\lim_{N \rightarrow \infty} \sigma(N) = 0, \quad \lim_{N \rightarrow \infty} N \sigma^N(N) = \infty$$

## RBF-Netze

Typisch: 2-Schichten Netzwerk



Aktivität

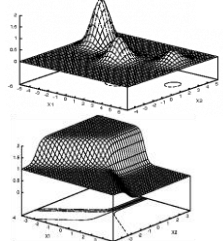
nicht normiert

$$f_i(x) = \sum_{k=1}^m w_k y_k = \sum_{k=1}^m w_k S_k(x)$$

mit  $S_k(x) = e^{-\frac{\|c_k - x\|^2}{2\sigma^2}}$

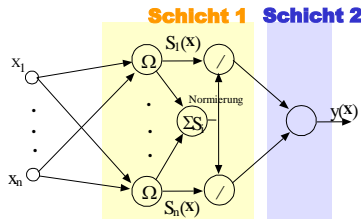
normiert

$$f_i(x) = \sum_{k=1}^m w_k y_k = \frac{\sum_{k=1}^m w_k S_k(x)}{\sum_{j=1}^m S_j(x)}$$



## RBF-Netze

### Aktivität Normiertes RBF-Netzwerk



$$y(x) = f(x) = \sum_i w_i \tilde{S}_i(x, c_i) \quad \text{mit} \quad \tilde{S}_i(x, c_i) = \frac{S_i(x, c_i)}{\sum_k S_k(x, c_k)}$$

Rüdiger Brause: Adaptive Systeme, Institut für Informatik, WS 2013/14

- 61 -

---

---

---

---

---

---

---

---

---

---

## Klassifikation mit RBF-Netzen

### Beste Klassifizierung *Bayes-Klassifizierung*

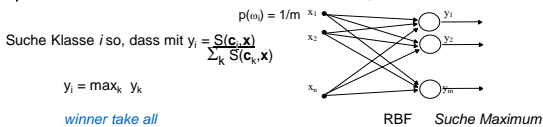
Suche Klasse  $\omega_i$  so, daß  $p(\omega_i|x) = \max_k p(\omega_k|x)$

$$\text{Wir wissen: } p(\omega_i|x) = \frac{p(\omega_i, x)}{p(x)} = \frac{p(\omega_i, x)}{\sum_k p(\omega_k, x)} = \frac{p(x|\omega_i)p(\omega_i)}{\sum_k p(x|\omega_k)p(\omega_k)}$$

Situation:

$$\text{Gaußvert. Abweichg. vom Klassenprototypen } c_i: p(x|c_i) = A e^{-\frac{(c_k - x)^2}{2\sigma^2}} = S(c_i, x)$$

➔ **Bayes-Klassifizierung mit NN:** Seien alle Klassen gleichwahrscheinlich



$$y_i = \max_k y_k$$

winner take all

R. Brause, Adaptive Modellierung: Kap. 2 Black-Box-Modellierung

- 242 -

---

---

---

---

---

---

---

---

---

---

## Klassifikation mit winner-take-all

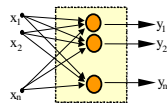
### Suche Maximum der Aktivität

#### Ein-Schicht-Netzwerk

(1 Cluster = 1 Klasse)

Suche Klasse  $i$  so, dass mit  $y_i = S(c_i, x) / \sum_k S(c_k, x)$

$$y_i = \max_k y_k$$

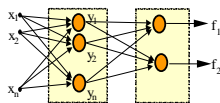


#### Zwei-Schichten-Netzwerk

(mehrere Cluster = 1 Klasse)

Suche Klasse  $i$  so, dass mit  $f_i = \sum_j w_j y_j$

$$f_i = \max_k f_k$$



⇒ Lernen **nur** der Gewichte für  $y_i$  bzw.  $f_i$

R. Brause, Adaptive Modellierung: Kap. 2 Black-Box-Modellierung

- 243 -

---

---

---

---

---

---

---

---

---

---

## Frage

### Was ist das Ziel der Bayes-Klassifikation?

#### Antwort

1. Die minimale Fehlerwahrscheinlichkeit
2. Die maximale bedingte Wahrscheinlichkeit für eine Entscheidung
3. Die minimale Abweichung vom korrekten Wert
4. Die maximale Korrelation mit dem korrekten Wert

## RBF- Lernen

### Lernmethode: Einzelne Anpassung der Schichten

- Anpassen der 1. Schicht: Lage  $c_i$  und Varianz  $C$
- Anpassen der 2. Schicht: Gewichte  $w_j$

**Pro:** schneller, da weniger komplex

**Contra:** Suboptima möglich

### Lernmethode: Gleichzeitige Anpassung beider Schichten

z.B. mit Backpropagation

**Pro:** Globales Optimum leichter erreichbar

**Contra:** langsam, mehr Daten für gegebenen max. Fehler nötig

## RBF- Lernen: 1.Schicht

### Erstellen der 1. Schicht

#### Bekannte Trainingsdaten

- Festlegen der RBF- Zentren und Weiten durch Clustering (k-mean, Kohonen-Netze, ...)
- Initiale Festlegung der **Anzahl** der Zentren, Iterative Verbesserung von Lage und Weite durch sequentielles Training

#### Unbekannte Daten

- Feste Unterteilung des Eingaberaumes, z.B. Netz potentieller Zentren bei uniformer, fester Weite für Datenbereich
- Inkrementeller Aufbau des Netzes: Sequentielle Regression

## RBF-Lernen : 1.Schicht

### • **k-means-Clustering**

Wähle  $k$  zufällige Muster  $x_j \in A$  als Clusterzentren  $c_j$ , bilde  $C_j = \{c_j\}$

#### REPEAT

- Ordne alle  $x_i \in A$  zu den nächstgelegenen Clusterzentren zu:  
Suche für  $x_i$  das Cluster  $c_z$  so, dass  $|x_i - c_z| = \min_k |x_i - c_k|$ , und füge  $x_i$  zu  $C_z$  zu.
- Entferne alle Cluster  $i$  mit  $|C_i| \leq 1$
- Bilde für jedes Cluster  $k$  ein neues Zentrum  $c_k = \langle x \rangle$  als Mittelwert aller Muster in  $C_k$

UNTIL Iterationszahl > Max

---

---

---

---

---

---

---

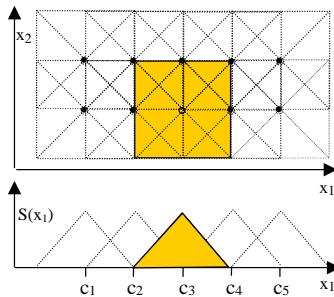
---

---

---

## RBF- Lernen : 1.Schicht

### • **Feste Unterteilung**




---

---

---

---

---

---

---

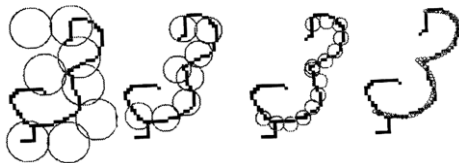
---

---

---

## RBF- Lernen : 1.Schicht

### • **Adaptive Unterteilung**




---

---

---

---

---

---

---

---

---

---

## RBF-Lernen

### Erstellen der 1. Schicht: Sequentielle Regression

- Start mit einem Neuron
- Füge ein neues Neuron hinzu für jedes Beispiel mit hohem Fehler (Abweichung vom gewünschten Netz-Ausgabewert)
- Verändere die Parameter bei den Nachbarn so, dass der Fehler verringert wird (Einpassen des neuen Neurons)

Das Netzwerk wächst solange,  
bis der Approximationsfehler auf das  
gewünschte Maß zurückgegangen ist.

## RBF-Lernen

### 2. Schicht: Fehlerminimierung MSE der lin. Schicht

z.B. mit Backpropagation-Lernen 2. Schicht

$$\Delta \mathbf{w}_i = -\gamma (y_i - L_i) \mathbf{x} \quad \text{Gewichte von } S_k \text{ zu Ausgabe } i$$

$$\mathbf{x} := (S_1, \dots, S_n) \text{ und } \mathbf{w} := (w_1, \dots, w_n)$$

oder: TLMSE

$$\mathbf{w}_i(t+1) = \mathbf{w}_i(t) - \gamma_t \mathbf{x} y_i \quad \text{und} \quad |\mathbf{w}_i(t+1)| = 1 \quad \text{Anti-Hebb Lernen}$$

$$\mathbf{x} := (S_1, \dots, S_n, L_i) \text{ und } \mathbf{w} := (w_1, \dots, w_n, w_{n+1})$$

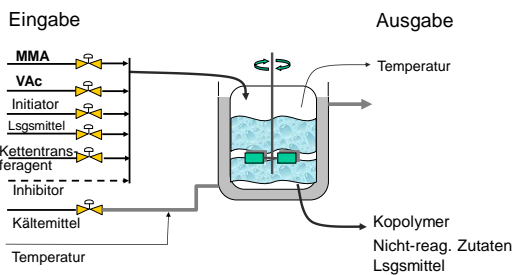
oder: Konkurrenz-Lernen

$$y_c = \max_i y_i$$

$$\Delta \mathbf{w}_c = -\gamma (y_c - L_c) \mathbf{x} \quad \text{Gewichte von } S_k \text{ zu Ausgabe } c$$

## Anwendung: Industriereaktor

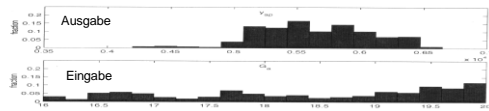
### Chem. Synthese eines Polymers



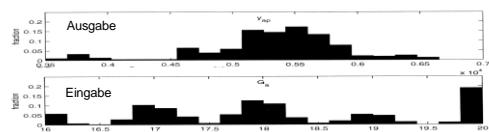


## Inputgenerierung

### Strategie **RUS**: Random uniform sampling



### Strategie **RDS**: Random distributed sampling

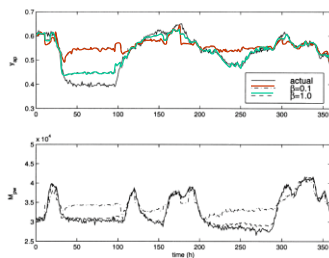


R.Brause, Adaptive Modellierung: Kap.2 Black-Box-Modellierung

- 273 -

## Ergebnisse

### Schmale vs. breite RBF

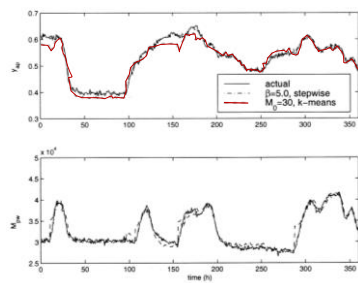


R.Brause, Adaptive Modellierung: Kap.2 Black-Box-Modellierung

- 274 -

## Ergebnisse

### K-means Clustering vs. Sequ. Regression



R.Brause, Adaptive Modellierung: Kap.2 Black-Box-Modellierung

- 275 -

## Black-Box- Modellierung

### Lineare Modellierung

### Nichtlin. Modellierung: MLP

### Nichtlin. Modellierung: RBF

## Eigenschaften von NN

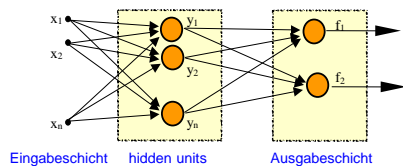
## Approximationsleistungen Neuronaler Netze

### • Allgemeine Eigenschaften von NN

- Kann man jede beliebige Funktion approximieren?
- Wenn ja, mit welcher Architektur ?
- Wie viele Schichten benötigt man ?
- Wie viele Neuronen pro Schicht braucht man ?

## Fähigkeiten: Nomenklatur

### • Mehrschichtennetze



$$\hat{f}(\mathbf{x}) = \sum_{j=1}^m w_j^{(2)} S(\mathbf{w}_j, \mathbf{x})$$

## Approximationsleistungen Neuronaler Netze

### Voraussetzungen

- **Sigma-Funktionen** (2-Schicht. NN)

$$\Sigma^n(S) := \{ \hat{f} \mid \hat{f} : \mathbb{R}^n \rightarrow \mathbb{R} \text{ mit } \hat{f}(\mathbf{x}) = \sum_{j=1}^m w_j^{(2)} S(z_j(\mathbf{x})) \}$$

wobei  $w_j^{(2)}$  aus  $\mathbb{R}$ ,  $\mathbf{x}$  aus  $\mathbb{R}^n$ ,  
und  $z_j \in \mathcal{Z}^n := \{ z \mid z(\mathbf{x}) = w^{(1)T} \mathbf{x} + b \}$  *affine Funktionen*  $\mathbb{R}^n \rightarrow \mathbb{R}$   
 $\lim_{x \rightarrow \infty} S(x) = 1, \quad \lim_{x \rightarrow -\infty} S(x) = 0$

## Approximationsleistungen Neuronaler Netze

### Aussage 1 (diskrete Punkte) *Hornik, Stinchcombe, White 1989*

Für die Funktionswerte *jeder beliebigen Funktion*

$$f(x) : \mathbb{R}^n \rightarrow \mathbb{R} \text{ von } N \text{ Mustern } x^1 \dots x^N$$

gibt es eine Sigma-Funktion  $f$ , so dass für alle Muster  $x^i$

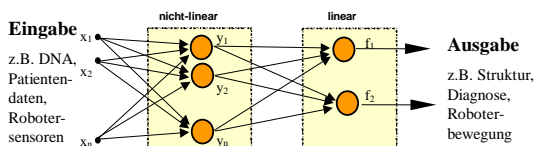
mit  $i = 1 \dots N$  gilt

$$\hat{f}(x^i) = f(x^i) \text{ punktweise Übereinstimmung}$$

## Approximationsleistung Neuronaler Netze

### Aussage 2 (stetige Funktionen)

Ein 2-Schichtennetzwerk mit nicht-linearer Ausgabefunktion  $S(z)$  kann *JEDER beliebige Funktion* beliebig dicht approximieren, wenn genügend Neuronen ex.



## Fähigkeiten der Multilayer-Netzwerke

### 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  $L_\infty$ -Norm in der Menge  $C^0$  aller stetigen Funktionen und  $\rho_0$ -*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