

# 6. Netze radialer Basisfunktionen

- 6.1 Beiträge zur Konzeptbildung
- 6.2 Definition von RBF-Netzen
- 6.3 Zweischichtig, sequentielles Lernen
- 6.4 Weitere Lernvariante für RBF-Netze
- 6.5 Approximation und Regularisierung
- 6.6 Statistische Interpretation von RBF-Netzen

# 6.1 Beiträge zur Konzeptbildung

Es sei eine Funktionsapproximation gewünscht.

Überblick:

- Satz von Kolmogorov (1957)
- Interpolation nach Powell (1987)
- Neuronale Informationsverarbeitung

# Satz von Kolmogorov (1957)

Gegeben sei eine  $I$ -dimensionale stetige Funktion  $f$ .

Dann existieren ein-dimensionale stetige Funktionen  $g_0$  und  $g_j$ ,  
 $j \in \{1, \dots, (2I + 1)\}$ , sowie Konstanten  $\alpha_i$ ,  $i \in \{1, \dots, I\}$ ,  
so daß gilt:

$$f(x_1, \dots, x_I) \stackrel{!}{=} \sum_{j=1}^{2I+1} g_0 \left( \sum_{i=1}^I \alpha_i g_j(x_i) \right)$$

# Satz von Kolmogorov (1957)

*Deutung:* Jede multi-dimensionale stetige Funktion kann durch Kombination von ein-dimensionalen stetigen Funktionen repräsentiert werden.

*Problem:* Satz ist nicht konstruktiv.

Es ist unbekannt, wie  $g_0$  und  $g_j$ ,  $j \in \{1, \dots, (2I + 1)\}$ , gefunden werden.

# Interpolation nach Powell (1987)

Gegeben seien  $M$  unterschiedliche Punkte

$$\{x^m \in \mathbb{R}^I | m = 1, \dots, M\}$$

und reelle Zahlen  $\{r^m \in \mathbb{R} | m = 1, \dots, M\}$ .

Dann findet man eine Funktion  $f : \mathbb{R}^I \rightarrow \mathbb{R}$  gemäß nachfolgender Definition, so daß die Interpolationsbedingung erfüllt ist:

$$f(x^m) := \sum_{j=1}^M w_j h_j(x^m) = r^m; \quad m \in \{1, \dots, M\}$$

# Interpolation nach Powell (1987)

Als Interpolationsfunktionen  $h_j$  unterscheidet man lokalisierende und nicht-lokalisierende Basisfunktionen  $h$ , die nicht-linear sind.

Dabei sei definiert:  $h_j(x) := h(\|x - x^j\|) =: h(d)$

Lokalisierende Basisfunktionen:

$$d \rightarrow \infty \Rightarrow h(d) \rightarrow 0$$

Beispiele:

$$h(d) := e^{-\left(\frac{d}{\sigma}\right)^2}; \quad \text{Gauß-Funktion}$$

$$h(d) := \frac{1}{(\gamma^2 + d^2)^\zeta}; \quad \zeta > 0$$

# Interpolation nach Powell (1987)

Nicht-lokalisierende Basisfunktionen:

$$d \rightarrow \infty \Rightarrow h(d) \rightarrow \infty \text{ oder} \\ h(d) \rightarrow \text{Konstante } (\neq 0)$$

Beispiele:

$$h(d) := (\gamma^2 + d^2)^\zeta; \quad \zeta > 0, \quad h(d) \rightarrow \infty$$

$$h(d) := \frac{1}{1+e^{-d}}; \quad h(d) \rightarrow 1$$

# Interpolation nach Powell (1987)

Bestimmung der Koeffizienten  $w_j$  zur Interpolation.

Angenommen, Basisfunktion  $h$  und zugehörige Matrix  $\mathbf{H}$  seien wie folgt gegeben:

$$\mathbf{H} := \begin{pmatrix} h_{11} & \cdots & h_{1M} \\ \vdots & & \vdots \\ h_{M1} & \cdots & h_{MM} \end{pmatrix}$$

$$h_{mj} := h(\|x^m - x^j\|)$$

$$m \in \{1, \dots, M\}, \quad j \in \{1, \dots, M\}$$



# Interpolation nach Powell (1987)

Weiterhin seien die Sollwerte gegeben:

$$\boldsymbol{R} := \begin{pmatrix} \boldsymbol{r}^1 \\ \vdots \\ \boldsymbol{r}^M \end{pmatrix}$$

Gesucht ist Vektor der Koeffizienten

$$\boldsymbol{w} := \begin{pmatrix} w_1 \\ \vdots \\ w_M \end{pmatrix}$$

$$\text{Lösung: } \boldsymbol{R} = \boldsymbol{H} \cdot \boldsymbol{w} \implies \boldsymbol{w}^* := \boldsymbol{H}^{-1} \cdot \boldsymbol{R}$$

# Neuronale Informationsverarbeitung

- Lokale Informationsverarbeitung  
(Neuronen sind sensitiv für lokalisierte Teilräume des Eingaberaumes mit Überlappungsgebieten)
- Partitionierung des Lernproblems

## 6.2 Definition von RBF-Netzen

### Überblick:

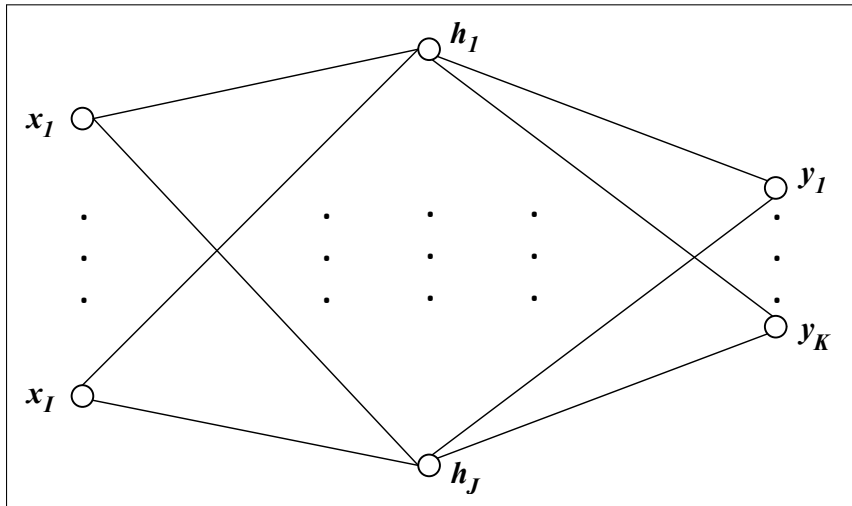
- Topologie und Basisfunktionen
- Funktion eines RBF-Netzes
- Funktionsapproximation durch RBF-Netz
- Parameter eines RBF-Netzes
- Nicht-isotrope Basisfunktionen
- Einschub: Zufallsvariable, Mittelwert, Varianz, Kovarianz, Kovarianzmatrix

# Topologie und Basisfunktionen

Ein RBF-Netz ist ein zweischichtiges Netz (d.h. eine verdeckte Schicht), deren Output-Knoten eine (gewichtete) Linearkombination von Projektionen auf radiale Basisfunktionen berechnen.

Die Basisfunktionen gehören zu den verdeckten Knoten, die eine lokalisierte Antwort auf den Input erzeugen.

# Topologie und Basisfunktionen



# Funktion eines RBF-Netzes

Lineare Kombination von nicht-linearen Basisfunktionen:

$$y_k(x) := \sum_{j=1}^J w_{jk} h_j(x); \quad k \in \{1, \dots, K\}$$

Wichtigste Basisfunktion:

Radial-symmetrischer Gauß (synonym Isotroper Gauß)

$$h_j(x) := h(\|x - \mu^j\|) := e^{-\frac{\|x - \mu^j\|^2}{2\sigma_j^2}}$$

# Funktionsapproximation durch RBF-Netz

RBF-Netz repräsentiert eine stetige Funktion, wobei die Anzahl der Basisfunktionen von der Komplexität der zu repräsentierenden Funktion abhängt.

Anzahl  $J$  der Basisfunktionen  $h_j$  ist klein im Vergleich zur Anzahl  $M$  der Trainingselemente.

⇒ Approximation statt Interpolation.

# Parameter eines RBF-Netzes

Zwei Parameterarten.

- Parameter der einzelnen RBFs:  
Zentren  $\mu^j$   
Distanzen  $\sigma_j$ , je zw. Zentrum und Wendepunkt einer Gauß-Kurve
- Parameter zur Kombination der RBFs:  
Faktoren (Gewichte)  $w_{jk}$



# Parameter eines RBF-Netzes

Basisfunktionen sind nicht mehr auf Trainingselementen lokalisiert.

⇒ Lokalisierung  $\mu^j$  ist Teil des Lernens.

Hinweis: Ergänzend zur Lokalisierung  $\mu^j$ , wird auch die Anzahl  $J$  der Basisfunktionen gelernt. D.h., die Anzahl der Knoten der verdeckten Schicht wird gelernt.

# Parameter eines RBF-Netzes

Jede Basisfunktion hat eine individuelle Einzugsweite, charakterisiert durch  $\sigma_j$ .

Falls  $\sigma_j$  klein/groß, dann ist die Basisfunktion für einen kleinen/großen Teilbereich des Eingaberaums zuständig.

Weiterhin gilt: Falls die  $\sigma_j$  der Basisfunktionen klein/groß, dann beeinflussen wenige/viele Basisfunktionen das Ergebnis der Netzanwendung auf ein bestimmtes Eingabelement.

⇒ Anpassung von  $\sigma_j$  ist Teil des Lernens.

# Parameter eines RBF-Netzes

Faktor  $w_{jk}$  gewichtet die Rolle der Basisfunktion  $h_j$  bei der Berechnung der  $k$ .ten Output-Komponente.

⇒ Ermittlung dieser Gewichte ist auch Teil des Lernens.

# Nicht-isotrope Basisfunktionen

Verallgemeinerung des radial-symmetrischen Gauß zu einem nicht-isotropen Gauß als Basisfunktion.

Definition der Basisfunktion mit Kovarianzmatrix  $S_j$ :

$$h_j(x) := e^{-\frac{1}{2}(x-\mu^j)^T(S_j)^{-1}(x-\mu^j)}$$

Im Falle eines  $I$ -dimensionalen Input-Vektors hat jede Basisfunktion dann  $I(I + 3)/2$  justierbare Parameter:  $I$  Parameter in  $\mu^j$  und  $I(I + 1)/2$  Parameter in  $S^j$ .

Balance erforderlich zwischen kleiner Zahl von (nicht-isotropen) Basisfunktionen mit vielen justierbaren Parametern und großer Zahl von (isotropen) Basisfunktionen mit wenigen Parametern.

# Einschub: Zufallsvariable, Mittelwert

Zufallsvariable  $X$

$M$  mögliche Realisierungen  $x^m$

z.B. beim Würfelspiel sechs mögliche Realisierungen  $1, \dots, 6$

Mittelwert  $\bar{X} := \frac{1}{M} \sum_{m=1}^M x^m$

# Einschub: Varianz, Kovarianz

Varianz  $\sigma_{xx} := \frac{1}{M} \sum_{m=1}^M (x^m - \bar{x})^2$

Zufallsvariablen  $X, Y$

Mittelwerte  $\bar{X}, \bar{Y}$

Kovarianz  $\sigma_{xy} := \frac{1}{M} \sum_{m=1}^M (x^m - \bar{x}) \cdot (y^m - \bar{y})$

# Einschub: Kovarianzmatrix

Zufallsvariablen

$X_1, \dots, X_I$

Kovarianz  $\sigma_{X_i X_j}$

Kovarianzmatrix  $S =$

$$\begin{pmatrix} \sigma_{X_1 X_1} & \dots & \sigma_{X_1 X_I} \\ \vdots & & \vdots \\ \sigma_{X_I X_1} & \dots & \sigma_{X_I X_I} \end{pmatrix}$$

$S$  ist symmetrisch, und  
hat  $I \cdot (I+1) / 2$  Parameter.

# Einschub: Kovarianzmatrix

Bsp.:



$$S = \begin{pmatrix} \sigma_{x_1 x_1} & \sigma_{x_1 x_2} \\ \sigma_{x_2 x_1} & \sigma_{x_2 x_2} \end{pmatrix}$$

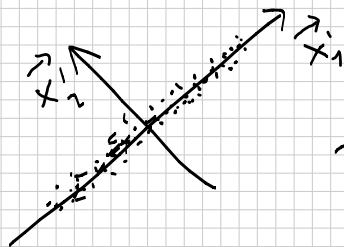
The matrix is annotated with circles and non-zero symbols:

- A circle around  $\sigma_{x_1 x_2}$  with  $\neq 0$  written above it.
- A circle around  $\sigma_{x_2 x_1}$  with  $\neq 0$  written below it.



# Einschub: Kovarianzmatrix

Bsp.:



$\leadsto$

$$S = \begin{pmatrix} \sigma_{x_1' x_1'} & 0 \\ 0 & \sigma_{x_2' x_2'} \end{pmatrix}$$

## 6.3 Zweischichtig, sequentielles Lernen

### Überblick:

- Annahmen für das Lernen
- Art des Lernens in den Schichten
- MEANS Cluster-Algorithmus
- ISODATA Cluster-Algorithmus
- Lineare Regression durch Pseudo-Inverse

# Annahmen für das Lernen

Es sei eine Funktionsapproximation gewünscht.

Seien  $(\boldsymbol{x}^m, r^m) \in \mathbb{R}^I \times \mathbb{R}$  die Trainingselemente, mit  $m \in \{1, \dots, M\}$ .

Seien  $h_j(\boldsymbol{x})$  die Symbole für die Basisfunktionen, mit  $j \in \{1, \dots, J\}$ .

Sei  $\boldsymbol{f}(\boldsymbol{x})$  das Symbol für den reellen Netzwerk-Output (oBdA skalar).

# Art des Lernens in den Schichten

1. Schicht: Unüberwachtes Lernen der Basisfunktionen bzgl. Eingabedaten  $x^m$  (keine Berücksichtigung von  $r^m$ ).  
Gelernt werden  $\mu^j, \sigma_j, J$ .
2. Schicht: Überwachtes Lernen der Faktoren  $w_j$  für die Linearkombination der Basisfunktionen  $h_j$ , durch Lineare Regression unter Berücksichtigung von  $r^m$ .

# MEANS Cluster-Algorithmus

Ziel: Unüberwachtes Lernen, d.h. Clustern von Trainingsdaten zur Bestimmung von Anzahl, Lage, Ausdehnung der Basisfunktionen.

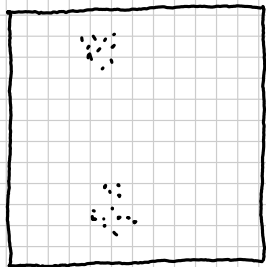
Methode: Partitionierung von  $\Omega_T$  in  $J$  disjunkte Teilmengen  $Cl_j$ , durch Minimierung der Intra-Cluster-Varianz  $V_{IC}$ .

$$V_{IC} := \frac{1}{J} \sum_{j=1}^J \left( \frac{1}{M_j} \sum_{x^m \in Cl_j} \|x^m - \mu^j\|^2 \right)$$

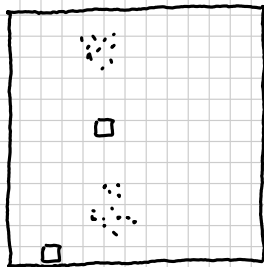
$$\mu^j := \frac{1}{M_j} \sum_{x^m \in Cl_j} x^m, \quad M_j := |Cl_j|$$

# MEANS Cluster-Algorithmus

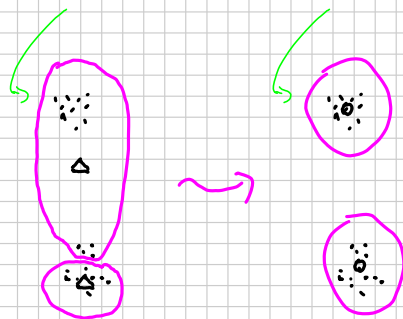
Bild,  
detektierte  
Punkte



Beliebige  
Initialisierung  
Cluster-Zentren



Intra-  
Cluster-Varianz  
groß klein



# MEANS Cluster-Algorithmus

Algorithmus im Überblick:

- Initialisierung der Zentren der Mengen  $Cl_j$ .
- Zugehörigkeiten von Elementen  $x^m$  zu den Mengen  $Cl_j$  festlegen bzw. ändern.
- Berechnung bzw. Neuberechnung der Zentren  $\mu^j$ .

# MEANS Cluster-Algorithmus

```
proc MEANS
{ Für alle j initialisiere  $\mu[j]$  beliebig
  Wiederhole solange sich  $Cl[j]$  ändern
  {
    Für alle  $x[m]$  aus  $\Omega[T]$  // Gruppierung
    {  $j\_b = \arg \min \{ ||x[m] - \mu[j]|| \}$ 
      füge  $x[m]$  zu  $Cl[j\_b]$  hinzu
    }
    Für alle  $Cl[j]$ 
    {  $M[j] = \text{Anzahl}(Cl[j])$  // Anzahl Elemente
       $\mu[j] = 1/M[j] * \text{Summe}(Cl[j])$  // Zentrum
    }
  }
}
```

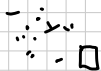


# MEANS Cluster-Algorithmus

## Bemerkungen:

- Bei jeder Iteration wird garantiert, daß  $V_{IC}$  nicht wächst.
- Manche Cluster bleiben eventuell leer.
- Resultierende Cluster hängen leider von Initialisierung ab.
- Zentren  $\mu_j$  der Cluster als Zentren der Basisfunktionen verwenden.
- Die Kovarianzmatrizen  $S_j$  der Cluster zur Definition von nicht-isotropen Basisfunktionen verwenden. Oder die Varianzen  $\sigma_j$  der Cluster zur Definition von isotropen Basisfunktion verwenden.

# MEANS Cluster-Algorithmus



Nur 3 Cluster,  
das vierte ist "leer".

Allg. entstehen bei

(7-) Means also

$\leq 7$  Cluster.

# ISODATA Cluster-Algorithmus

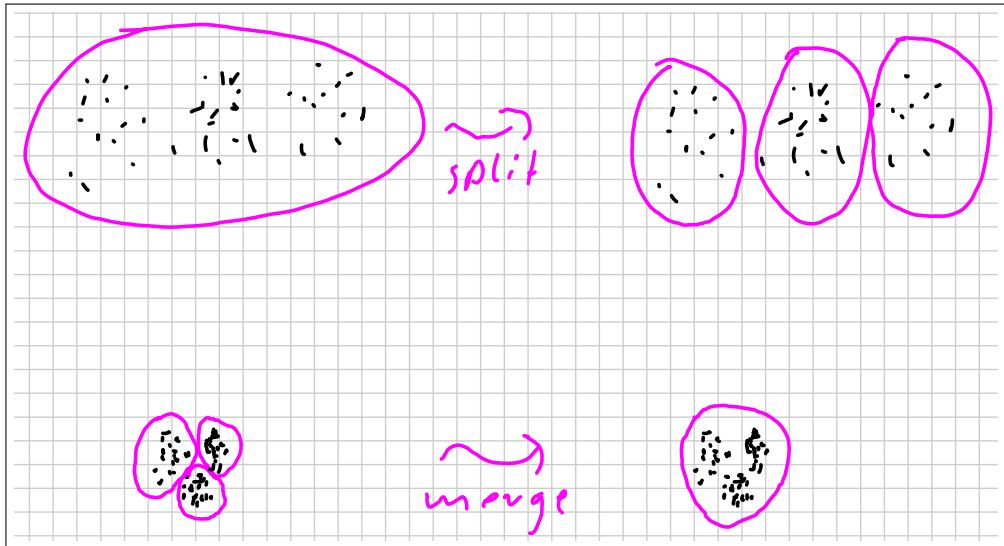
Iterative **S**elf-**O**rganizing **D**ATA Analysis Technique.

Algorithmus im Überblick:

- Die Anzahl der Cluster ist variabel.
- Start mit Vorgabe niedriger Zahl von Clustern.
- Variieren der Größe und Gestalt der Cluster.
- Anfangsclusterung mit MEANS.
- Split-Operation falls Streuung im Cluster zu groß.
- Merge-Operation falls die Zentren der Cluster zu eng benachbart.

Funktioniert besser als MEANS (reduzierte Abhängigkeit von Initialisierung).

# ISODATA Cluster-Algorithmus



# Lineare Regression durch Pseudo-Inverse

- Gewichtsvektoren zwischen verdeckter Schicht und Ausgabeschicht werden unter der Annahme gelernt, dass die Basisfunktionen der verdeckten Schicht vorliegen.
- Diese Basisfunktionen werden auf die Trainingselemente angewendet.
- Weil die Knoten der Ausgabeschicht als Propagierungsfunktion den Linearen Assoziator und als Aktivierungsfunktion die Identität haben, ist nur noch eine Lineare Regression erforderlich.
- Beim Batch-Lernen erfolgt die Lineare Regression mit der Pseudo-Inversen (siehe Unterkapitel 4.5), und führt zur Minimierung des Mean-Squared-Error.

# Lineare Regression durch Pseudo-Inverse

Funktion des RBF-Netzes:  $f(x) := \sum_{j=1}^J w_j h_j(x)$

Trainingsdaten:  $(x^m, r^m) \in \Omega_T \subset \mathbb{R}^I \times \mathbb{R}$

Fehlerfunktion:  $D(w) := \frac{1}{2} \sum_{m=1}^M (f(x^m) - r^m)^2$

Optimaler Gewichtsfaktor:  $w^* := \underbrace{(H^T \cdot H)^{-1} \cdot H^T \cdot R}_{\text{Pseudo-Inverse von } H}$

Es ist  $w^*$  der Vektor mit den  $J$  optimalen Gewichten,  
 $H$  die  $(M \times J)$ -Matrix mit Komponenten  $h_j(x^m)$ ,  
und  $R$  der Vektor mit den  $M$  Solldaten.

## 6.4 Weitere Lernvariante für RBF-Netze

### Überblick:

- Diskussion über das zweistufige Lernen
- Integriertes Lernen aller Parameter
- Zweistufiges und integriertes Lernen

# Diskussion über das zweistufige Lernen

## Vorteile:

- Partitionierung des Lernproblems durch zweistufiges Lernen von zwei Parameterarten. Dies bewirkt eine Effizienzsteigerung beim Lernen.
- Transformation der Eingabe-(Roh-)Daten in einen anderen Raum, in welchem das Problem linearer Natur ist.
- Lernen der Gewichtsvektoren (in der zweiten Stufe) durch einfache, lineare Regression.



# Diskussion über das zweistufige Lernen

## Nachteile:

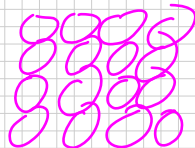
- Die unabhängige Optimierung der Basisfunktionen von Solldaten  $r^m$  kann ein Problem sein.
- Zur Abtastung eines  $I$ -dimensionalen Eingaberaumes wächst die notwendige Zahl der Basisfunktionen exponentiell mit  $I$ .
- Oft haben bestimmte Dimensionen in den mehr-dimensionalen Eingabedaten keinen Einfluß auf die Ausgabedaten. Dies führt aber nicht zur Verringerung der Zahl der Basisfunktionen.

# Diskussion über das zweistufige Lernen

Abtastung Eingabebereich

z.B. mit 4 Basisfunktionen je Dimension,

Somit  $4 \times 4$  bei 2 Dimensionen



Bzw.  $4^3$  bei 3 Dimensionen.

Exponentieller Zuwachs mit Dimensionszahl.

# Diskussion über das zweistufige Lernen

Belanglose Dimensionen



Belanglose Dimension  $\vec{x}_2$ , weil

$$\text{z.B. } f(x_{11}, x_{21}) = f(x_{11}, x_{22}) = f(x_{11}, x_{23}) = \dots$$

Dies gelte für alle Punkte auf  $\vec{x}_1$ -Achse.

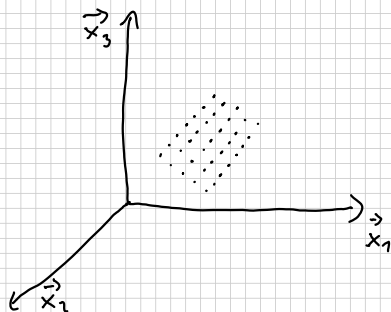
# Diskussion über das zweistufige Lernen

## Nachteile:

- Wenn die intrinsische Dimensionalität der Eingabedaten und die Lage des eingebetteten Unterraumes bekannt wäre, dann könnte man besser geeignete Basisfunktionen definieren, um deren Anzahl zu reduzieren.
- Die Anzahl der Basisfunktionen sollte in denjenigen Teilgebieten des Eingaberaumes dicht verteilt sein, wo der Abbildungsfehler groß ist, z.B. wo Funktionen stark gekrümmt sind.

# Diskussion über das zweistufige Lernen

Intrinsische Dimension



Bsp.: 3D - Eingaberaum, aber Trainingsmenge ist intrinsisch nur 2D.

# Diskussion über das zweistufige Lernen

Gewünschte Regressionsfunkt.



Ideale Verteilung der  
Basisfunktionen



# Integriertes Lernen aller Parameter

- Integriertes, überwachtes Lernen aller Parameter der zwei Parameterarten, d.h. keine Organisation in zwei Stufen, auch Basisfunktionen überwacht lernen.
- Gradientenabstiegsverfahren in nicht-konvexer Fehlerfunktion.
- Realisiert eine ganzheitliche, nicht-lineare Regression.
- Verfahren ermittelt eventuell nur lokales Minimum (statt globales) der Fehlerfunktion.
- Verfahren würde eine konstante Zahl von Basisfunktionen zu Grunde legen, also die Anzahl der verdeckten Knoten nicht lernen.

# Zweistufiges und integriertes Lernen

- Die Vorteile des zweistufigen Lernens und des integrierten Lernens könnten in folgendem Verfahren kombiniert werden.
- Erst zweistufiges Lernen, dann integriertes Lernen anschließen.
- Die Ergebnisse des zweistufigen Lernens dienen zur Initialisierung des Gradientenabstiegs beim integrierten Lernen.
- Dadurch wird Effizienz gesteigert und die Chance auf Erreichen des globalen Minimums der Fehlerfunktion erhöht.



# 6.5 Approximation und Regularisierung

## Überblick:

- Gut oder schlecht gestellte Probleme
- Beispiele für schlecht gestellte Probleme
- Regularisierung
- Lineare Regression und Regularisierung

# Gut oder schlecht gestellte Probleme

Bezugnahme auf Hadamard (1923).

Gut gestelltes (well-posed) Problem:

- Es existiert eine Lösung, und die
- Lösung ist eindeutig, und die
- Lösung hängt kontinuierlich von Eingabedaten ab.

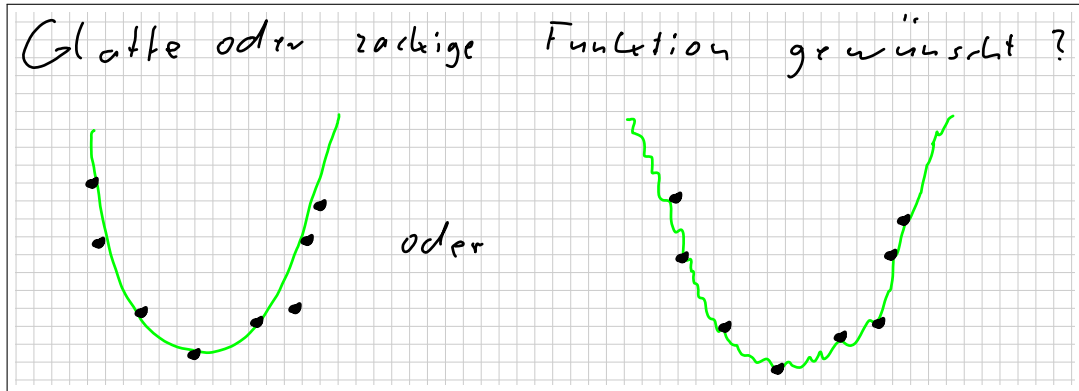
Ist eine dieser Bedingungen nicht erfüllt, so ist das Problem schlecht gestellt (ill-posed).

# Beispiele für schlecht gestellte Probleme

- Lineare Algebra:  $y = Ax \xRightarrow{?} x = A^{-1}y$
- Regelungstheorie: Prozeßidentifikation aus Beobachtung
- Computer Vision: Szenenrekonstruktion aus Bildern

Änderung/Verbesserung der Problemeigenschaften  
(ill-posed  $\Rightarrow$  well-posed) durch:

- mehr Meßergebnisse,
- Vorverarbeitung der Meßergebnisse,
- Definition von Zwängen durch zusätzliches Wissen.



# Regularisierung

Bezugnahme auf Poggio, Girosi (1990).

$$G(f) := \frac{1}{2} \sum_{m=1}^M (f(x^m) - r^m)^2 + \frac{\lambda}{2} \int |R_f|^2 dx$$

Generischer Regularisierungsoperator  $R_f$  für die Funktion  $f$ , der z.B. die Glattheit von  $f$  mißt.

Rolle (Wichtigkeit) des Regularisierungsterms wird festgelegt mit Hilfe des Regularisierungsparameters  $\lambda$ .

Bemerkung: Es gibt auch Regularisierungsterme mit mehreren Regularisierungsparametern.

# Lineare Regression und Regularisierung

Verwendung der empirischen Fehlerfunktion  $D(w)$  in der Definition einer allgemeineren Kostenfunktion (loss function)  $G(w)$ :

$$G(w) := D(w) + \underbrace{\frac{1}{2} \sum_{j=1}^J \lambda_j w_j^2}_{\text{Bestrafungsterm}}$$

Der zusätzliche Term bewirkt einen “Zwang” bezüglich der zu lernenden Funktion  $f$ , z.B. “bestraft” große Gewichte  $w_j$ .

Die Regularisierungsparameter  $\lambda_j \geq 0$  kontrollieren das Ausmaß der Bestrafung:

$\lambda_j$  klein  $\Rightarrow$  dichte Anpassung an Messungen

$\lambda_j$  groß  $\Rightarrow$  Verzicht auf dichte Anpassung

# Lineare Regression und Regularisierung

Durch "Bestrafungsterm" kann die Art der zu lernenden Funktion beeinflusst werden.



Verwendung 3 Gauße,  
|w| groß, dichte  
Anpassung an Daten



Verwendung 3 Gauße,  
|w| kleiner, grobe  
Anpassung an Daten.



# Lineare Regression und Regularisierung

$$\frac{\partial G}{\partial w_j} = \underbrace{\sum_{m=1}^M (f(x^m) - r^m)}_{\frac{\partial D}{\partial w_j}} \underbrace{\frac{\partial f}{\partial w_j}(x^m)}_{h_j(x^m)} + \lambda_j w_j \stackrel{!}{=} 0; \quad j \in \{1, \dots, J\}$$

$$\sum_{m=1}^M f(x^m) h_j(x^m) + \lambda_j w_j = \sum_{m=1}^M r^m h_j(x^m)$$

# Lineare Regression und Regularisierung

Vektor-Notation:

$$\vec{h}^{sp_j, T} \cdot F + \lambda_j w_j = \vec{h}^{sp_j, T} \cdot R$$

mit

$$\vec{h}^{sp_j} := (h_j(x^1), \dots, h_j(x^M))^T$$

$$F := (f(x^1), \dots, f(x^M))^T$$

$$R := (r^1, \dots, r^M)^T$$

Nun Stapel von  $J$  Gleichungen zusammenfassen.

# Lineare Regression und Regularisierung

Matrix-Notation:

$$H^T \cdot F + L \cdot w = H^T \cdot R$$

mit

$$\begin{aligned} H &:= \begin{pmatrix} h_1(x^1) & \dots & h_J(x^1) \\ \vdots & & \vdots \\ h_1(x^M) & \dots & h_J(x^M) \end{pmatrix} =: \begin{pmatrix} \vec{h}^{ze_1} \\ \vdots \\ \vec{h}^{ze_M} \end{pmatrix} \\ &=: (\vec{h}^{sp_1} \quad \dots \quad \vec{h}^{sp_J}) \end{aligned}$$

und

$$L := \begin{pmatrix} \lambda_1 & \dots & 0 \\ & \ddots & \\ 0 & \dots & \lambda_J \end{pmatrix}; \quad w := (w_1, \dots, w_J)^T$$

# Lineare Regression und Regularisierung

Es gilt:

$$f(x^m) = \sum_{j=1}^J w_j h_j(x^m) = \vec{h}^{ze_m} \cdot w$$

Zusammenfassung für alle  $m \in \{1, \dots, M\}$ :  $F = H \cdot w$

Verwendung in erster Formel von letzter Folie:

$$H^T \cdot R = H^T \cdot F + L \cdot w = H^T \cdot H \cdot w + L \cdot w = (H^T \cdot H + L) \cdot w$$

Ergibt optimalen Gewichtsvektor unter Zwang:

$$w^* := (H^T \cdot H + L)^{-1} \cdot H^T \cdot R$$

$\rightsquigarrow$  allg. Pseudoinverse

## 6.6 Statistische Interpretation von RBF-Netzen

### Überblick:

- Funktion von RBF-Netz als Bayes-Formel
- Probabilistische Deutung von RBF-Netzen

Annahme: RBF-Netze zur Klassifikation

# Funktion von RBF-Netz als Bayes-Formel

$$P(c^k|x) = \frac{P(x|c^k) \cdot P(c^k)}{P(x)} \quad (1)$$

$$P(x|c^k) = \sum_{j=1}^J P(x|j) \cdot P(j|c^k) \quad (2)$$

$$P(x) \stackrel{!}{=} \sum_{j'} P(x|j') \cdot P(j') \quad (3)$$

Deutung von  $P(x|j)$ : Zugehörigkeit von  $x$  zu einer bestimmten Gauß-Verteilung (Nummer  $j$ ).

Deutung von  $P(j|c^k)$ : Relevanz der Gauß-Funktion  $j$  für die Charakterisierung der Klasse  $c^k$ .

# Funktion von RBF-Netz als Bayes-Formel

$$\begin{aligned}P(x) &= \sum_{k=1}^K P(x|c^k) \cdot P(c^k) \\&= \sum_{k=1}^K \left( \sum_{j=1}^J P(x|j) \cdot P(j|c^k) \right) \cdot P(c^k) \\&= \sum_{j=1}^J \sum_{k=1}^K \underbrace{P(x|j)}_{\text{unabh. von } \sum_k} \cdot P(j|c^k) \cdot P(c^k) \\&= \sum_{j=1}^J P(x|j) \cdot P(j)\end{aligned}$$

# Funktion von RBF-Netz als Bayes-Formel

Einsetzen von Gleichung 2 und 3 in Gleichung 1:

$$\begin{aligned} P(c^k|x) &= \frac{\sum_{j=1}^J P(x|j) \cdot P(j|c^k) \cdot P(c^k) \cdot \frac{P(j)}{P(j)}}{\sum_{j'} P(x|j') \cdot P(j')} \\ &= \sum_{j=1}^J \underbrace{\frac{P(j|c^k) \cdot P(c^k)}{P(j)}}_{\substack{=P(c^k|j) \\ =w_{jk}}} \cdot \underbrace{\frac{P(x|j) \cdot P(j)}{\sum_{j'} P(x|j') \cdot P(j')}}_{\substack{=P(j|x) \\ =h_j(x)}} \end{aligned}$$



# Probabilistische Deutung von RBF-Netzen

Das Lernen eines RBF-Netzes wird so realisiert, dass sich für einen Input-Vektor  $x$  die a posteriori Wahrscheinlichkeit für eine Klasse  $c^k$  ergibt.

Eine Basisfunktion  $h_j$  gibt die Wahrscheinlichkeit der Zugehörigkeit des Input-Vektors  $x$  zur Gauß-Funktion mit Nummer  $j$  an.

Das Gewicht  $w_{jk}$  gibt die Relevanz der Gauß-Funktion  $j$  für die Charakterisierung der Klasse  $c^k$  an, ausgedrückt als bedingte Wahrscheinlichkeit.