

Radiale Basis- funktionen

AS2-5

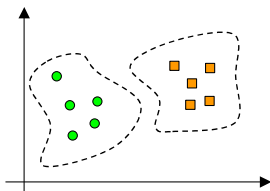
Approximation & Klassifikation mit RBF

Lernen in RBF-Netzen
support vector-Maschinen
Anwendung RBF-Netze

Radiale Basisfunktionen

Motivation: lokale Cluster-Klassenbildung

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



Radiale Basisfunktionen

Definition Glockenfunktionen

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 $a > 0$ mit $S_G(z)$ nicht anwachsend $\forall z \in [a, \infty)$,
nicht abfallend $\forall z \in (-\infty, a)$

Also ist $S_G(a)$ globales Maximum.

Glockenfunktionen

Beispiele

- **Kombination** von Quetschfunktionen

$$S_G(x_1, \dots, x_n) = \max(0, 1 - \sum_{i=1}^n b(x_i) - 1) \text{ mit } b(x) = \frac{S_Q(1+x_1) + S_Q(1-x_1) - 1}{2S_Q(0) - 1}$$
- **Ableitungen** von Quetschfunktionen

$$S_G(x) = \frac{\partial S_Q}{\partial x}$$
- **Produkte** von Glockenfunktionen

$$S_G(x_1, \dots, x_n) = S_G(x_1) \times \dots \times S_G(x_n)$$
- **allgemeine Radiale Basisfunktionen**

$$S_G(x) = h(|x|), \quad x \in \mathbb{R}^n, \quad h(\cdot) \text{ streng monoton fallend}$$
- aus Intervallen **zusammengesetzte** Funktionen

$$S_G(z) = (1-z^2)^{20} \text{ im Intervall } z \in [-1, +1], \text{ sonst null.}$$

RBF-Netze

Typisch: 2-Schichten Netzwerk

Aktivität

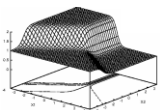
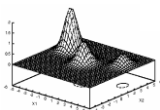
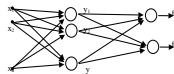
nicht normiert

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

mit $S_k(\mathbf{c}_k, \mathbf{x}) = e^{-\frac{(\mathbf{c}_k - \mathbf{x})^2}{2\sigma^2}}$

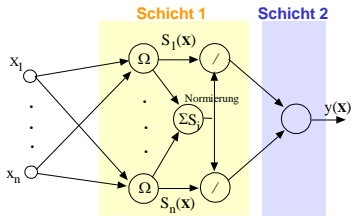
normiert

$$\hat{f}_i(\mathbf{x}) = \sum_{k=1}^m w_k y_k = \frac{\sum_{k=1}^m w_k S_k(\mathbf{x})}{\sum_{j=1}^m S_j(\mathbf{x})}$$



Radiale Basisfunktionen

Aktivität Normiertes RBF-Netzwerk



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

Radiale Basisfunktionen

Basisfunktionen maximaler Information (Entropie)

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

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**

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

$$\frac{\partial L}{\partial p} = 0, \quad \frac{\partial L}{\partial \mu_1} = 0 \quad (\text{Rechnung Kap.5.2})$$

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

Radiale Basisfunktionen

Basisfunktionen maximaler Information (Entropie)

$$H(p^*) = \max_p H(p) \quad x \in [0, 1], \quad p^*(x) = ?$$

NB: $\int_0^1 p(x) dx = 1$ oder $g(x) := \int_0^1 p(x) dx - 1 = 0$ *ausreichende NB*

Ansatz **Lagrange-Funktion**

$$L(p, \mu) := H(p) + \mu g(p)$$

$$\frac{\partial L}{\partial p} = 0 = \frac{\partial L}{\partial \mu}, \quad \frac{\partial L}{\partial \mu} = 0 \quad (\text{Rechnung analog Kap.5.2})$$

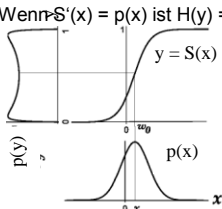
Ergebnis: $p^*(x) = \text{const}$ **Uniforme Verteilung**

Transformation mit maximaler Information

$[-\infty, +\infty] \ni x \rightarrow [0, 1]$ **Max. Information** bei uniformer pdf !

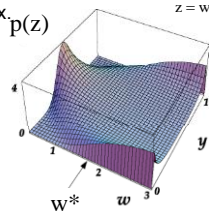
Wie ? (Rechnung Anhang A.4)

Wenn $S'(x) = p(x)$ ist $H(y) = \max_x p(z)$



$$S_p(z) = \frac{1}{1 + \exp(-z)}$$

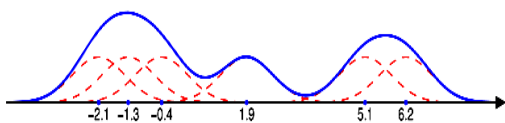
$$z = wx + w_0$$



Einstellung von $S(x)$ mittels w

Parzen Window - Methode

Approximation durch Überlagerung von Basisfunktionen



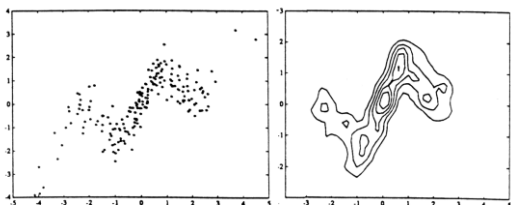
$$\rho(x) = \frac{1}{N} \sum_{i=1}^N W(x - x_i) \quad W(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-x^2/2\sigma^2}$$

Perfekte Approximation bei abnehmender Breite σ , wobei

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

Parzen Window

Approximation durch Überlagerung von Basisfunktionen



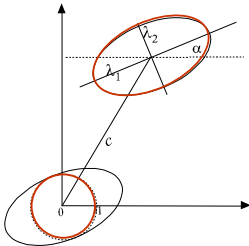
Normierung der Variablen

Problem

- PCA etc. problematisch bei heterogenen Variablen, z.B. (x_1 [cm], x_2 [Pascal], x_3 [°C])
Welche Einheiten pro Dimension?
- Welche Relation sollen die Einheiten zueinander haben ?

Normierung der Variablen

Lösung einheitliche Transformation aller Variablen durch Skalierung **S**, Drehung **D**, Verschiebung **V**



$$\mathbf{x} \rightarrow \mathbf{z} = \mathbf{SDV}\mathbf{x} = \mathbf{M}\mathbf{x}$$

$$d^2 = \mathbf{z}^T \mathbf{z} = \mathbf{x}^T \mathbf{M}^T \mathbf{M} \mathbf{x}$$

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

Mahalanobis-Abstand

RBF-Ausgabefunktion

$$S_G(\mathbf{x}) = A \exp(-(\mathbf{x} - \mathbf{c})^T \mathbf{C}^{-1} (\mathbf{x} - \mathbf{c}))$$

Klassifikation mit RBF-Netzen

Beste Klassifizierung

Suche Klasse ω_i so, daß $p(\omega_i|\mathbf{x}) = \max_j p(\omega_j|\mathbf{x})$ *Bayes-Klassifizierung*

Wir wissen: $p(\omega_i|\mathbf{x}) = \frac{p(\omega_i, \mathbf{x})}{p(\mathbf{x})} = \frac{p(\omega_i, \mathbf{x})}{\sum_j p(\omega_j, \mathbf{x})}$

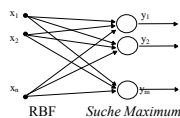
Annahme: Gaußverteilte Abweichungen der \mathbf{x} von den Klassenprototypen \mathbf{c}_j ,

also $p(\mathbf{c}_j, \mathbf{x}) = A e^{-\frac{(\mathbf{c}_j - \mathbf{x})^2}{2\sigma^2}} =: S(\mathbf{c}_j, \mathbf{x})$

⇒ *Bayes-Klassifizierung mit NN:*

Suche Klasse ω_k so, daß mit $Y_i = \frac{S_i(\mathbf{c}_i, \mathbf{x})}{\sum_j S_j(\mathbf{c}_j, \mathbf{x})}$

$y_k = \max_i y_i$ *winner take all*



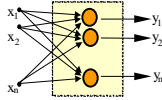
Klassifikation mit winner-take-all

Suche Maximum der Aktivität

Ein-Schicht-Netzwerk

Suche Klasse k so, dass mit $y_i = S(c_i, x) / \sum_j S(c_j, x)$

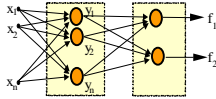
$$y_k = \max_i y_i$$



Zwei-Schichten-Netzwerk

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

$$f_i = \max_k f_k$$



⇒ Lernen **nur** der Gewichte für y_j bzw. f_i

Approximation & Klassifikation mit RBF

Lernen in RBF-Netzen

support vector-Maschinen

Anwendung RBF-Netze

Lernverfahren

Ansätze

- Schichtweise Einzelanpassung
 - Anpassen der ersten Schicht (Zentrum +Breite)
 - Anpassen der zweiten Schicht (Gewichte)
- Gesamtanpassung, z.B. durch Backpropagation

Anpassung der ersten Schicht

Phasen

1. **initiale Verteilung** (Anzahl, Lage und Form) der Glockenfunktionen
2. **iterative Adaption** der RBF-Parameter an die Trainingsdaten

Initiale Verteilung

- **Bekannte Trainingsdaten**
Clustersuche, RBF-Zentren = Clusterzentren; RBF-Breite = Clusterstreuung
- **Unbekannte Trainingsdaten**
 - Sukzessiver Netzaufbau
 - Überdeckung durch Fehlerminimierung
 - Überdeckung durch regelmäßiges Raster
 - Clusteranalyse durch Kohonen-Netze

Anpassung der ersten Schicht

Initiale Verteilung

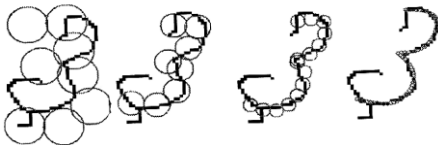
Sukzessiver, fehlerorientierter Netzaufbau

- 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.

Anpassung der ersten Schicht

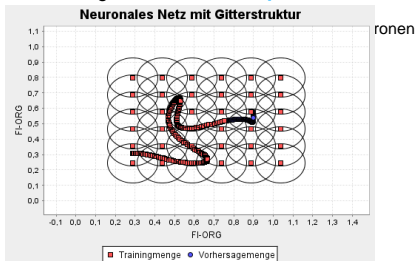
Initiale Verteilung

Adaptiver und sukzessiver Netzaufbau für Abdeckung einer Testverteilung



RBF-Probleme

- Sigmoidale Ausgabeffkt auch für Extrapolation,
- RBF-Ausgabeffkt nur für Intrapolation.

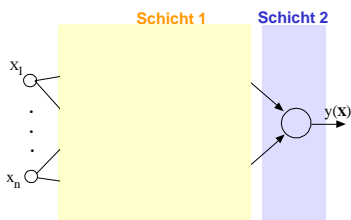


Rüdiger Brause: Adaptive Systeme, Institut für Informatik, WS 2013/14

- 22 -

Anpassung der zweiten Schicht

Normiertes RBF-Netz



$$y(\mathbf{x}) = \hat{f}(\mathbf{x}) = \sum_i w_i v_i \quad \text{mit } v_i = \tilde{S}_i(\mathbf{x}, c_i)$$

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

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

- 23 -

Anpassung der zweiten Schicht

• TLMSE: Eigenvektor fitting

$$\mathbf{w}(t) = \mathbf{w}(t-1) - \gamma(t) \mathbf{y} [\mathbf{x}(t) - \mathbf{w}(t-1)\mathbf{y}] \quad \text{negative Oja Lernregel}$$

mit Mittelwertskorrektur $\mathbf{y} = (\mathbf{x} - \mathbf{x}_0)^T \mathbf{w}$

• Minimisierung der Entropie

$$\mathbf{w}(t) = \mathbf{w}(t-1) - \gamma \text{grad } H(\mathbf{y}(\mathbf{w}))$$

Approximation von $p(\mathbf{x})$ mit Parzen Windows:

Rechnung

$$\mathbf{w}(t) = \mathbf{w}(t-1) - \gamma(y - y_k) (\mathbf{x} - \mathbf{x}_k) \quad \text{Hebb'sche Regel}$$

Ausgabe y , frühere Ein/Ausgabe k

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

- 24 -

Approximation & Klassifikation mit RBF

Lernen in RBF-Netzen

support vector-Maschinen

Anwendung RBF-Netze

Gesamtanpassung

• Lernen mit Backpropagation

Zielfunktion $R(M) = \langle (f(x, M) - F(x))^2 \rangle = \langle r(x, M) \rangle$

1. Schicht: Lernen der RBF-Koeffizienten M_{ij} durch Gradientenalgorithmus

$$M_{ij}^k(t+1) = M_{ij}^k(t) - \gamma \frac{\partial}{\partial M_{ij}^k} r(x, M^k)$$

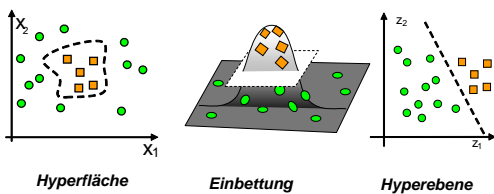
2. Schicht: Standard, z.B. BP

• Klassifikation durch support vector-Maschinen

Gesamtanpassung: nicht-lin. Separierung

• Klassifikation

Idee: Verwenden von RBF für lineare Separierung

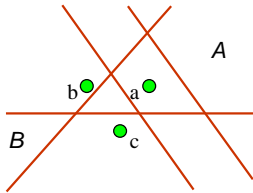


Gesamtanpassung

„Diagnosevariabilität“ h VC-Dimension

h = maximale Anzahl der Punkte, die durch die Diagnosemaschine auf 2^h Arten in zwei Klassen geteilt werden können.

Beispiel $h = 3$ Punkte, $2^3=8$ Diagnosearten möglich. $h < 4$



- $A=\{\}, B=\{a,b,c\}$
- $A=\{a\}, B=\{b,c\}$
- $A=\{a,b\}, B=\{c\}$
- $A=\{a,c\}, B=\{b\}$

sowie 4 Möglichkeiten bei Umbenennung $A \rightarrow B$, $B \rightarrow A$

Gesamtanpassung: support vector machine

Forderung für lin. Separierung

„Lege die Hyperebene so, dass sie maximalen Abstand zu allen Grenzpunkten hat“

$$\min_i |w^T z_i + b| = 1 \quad \text{Mindestabstand} = 1$$

$$f_{w,b}(z_i) = \text{sgn}(w^T z_i + b) = y_i \in \{+1, -1\} \quad \text{Klassifizierung}$$

$$(w^T z_i + b) y_i \geq 1$$

$$(w^T z_i + b) y_i \geq 1 - \xi_i \quad \begin{array}{l} \text{Minimierung des strukturellen Risikos} \\ \text{Schlupfvariable} \end{array}$$

Gesamtanpassung : support vector machine

Ansatz support vector – Maschine

- Alle Muster sind in einem Cluster: $|z_i - a| < r$ Kugelradius
- Endliche Beschreibung der Trennung $|w| \leq A$

$$\Rightarrow h < r^2 A^2 + 1 \quad \text{Vapnik 1995}$$

Reduzierung des Klassifizierungsfehlers durch Beschränkung von h

Neues Ziel: Minimierung von

$$T(w, \xi_i) = \frac{1}{2} w^2 + \gamma \sum_{i=1}^N \xi_i$$

$$\text{mit NB } g(w, i) = 1 - (w^T z_i + b) y_i - \xi_i = 0$$

support vector - Maschine

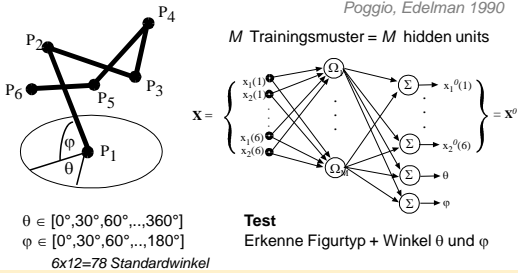
Ergebnisse

Daten	Klassifikations-Fehlerrate		
	Klass. RBF	RBF mit SV Zentren	Reine SV-Maschine
US Postal Service			
Training (7291 Muster)	1,7 %	0,0 %	0,0 %
Test (2007 Muster)	6,7 %	4,9 %	4,2 %

Frage: Warum ist diese Gegenüberstellung problematisch ?

Erkennen von 3D-Figuren

Training Feste Figur x_i aus 6 Punkten,
40-100 Random-Projekt. auf 2D-Fläche



Erkennen von 3D-Figuren

Ergebnisse Erkennungsleistung

