

# 8. Support-Vektor-Maschinen

Synonym: Support-Vektor-Netze.

Verwendung: Klassifikation oder Funktionsapproximation.

Hier: Fokussierung auf Klassifikation (zwei Klassen).

8.1 SVM mit linearer Diskriminanzfunktion

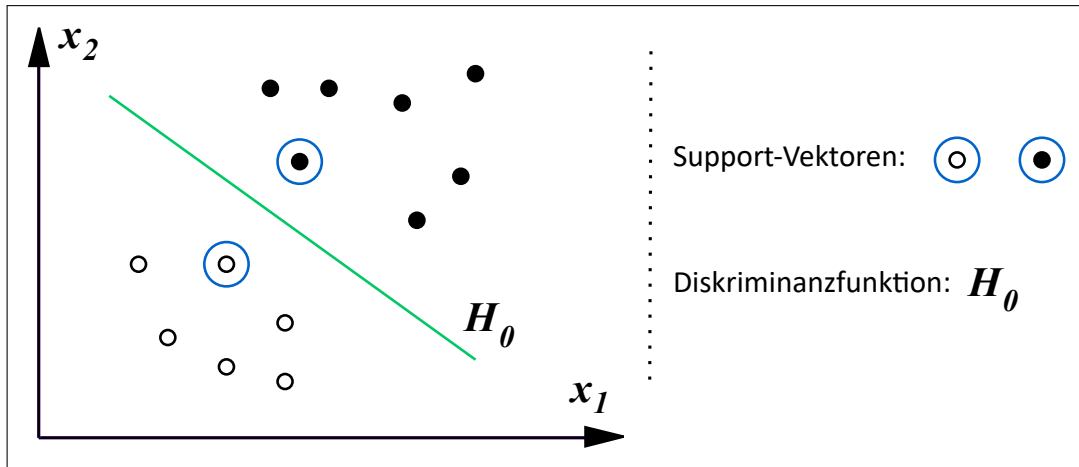
8.2 SVM mit nicht-linearer Diskriminanzfunktion

# 8.1 SVM mit linearer Diskriminanzfunktion

## Überblick:

- Support-Vektoren, Diskriminanzfunktion
- Trennband maximaler Breite
- Quadratische Optimierungsaufgabe
- Hinweis auf die optimale separierende Hyperebene
- Äquivalente quadratische Optimierungsaufgabe
- Berechnung der optimalen separierenden Hyperebene

# Support-Vektoren, Diskriminanzfunktion



# Trennband maximaler Breite

Aufgabe: Bestimmung einer Diskriminanzfunktion, so daß das dazu parallele Trennband maximale Breite hat.

Menge  $\Omega_T := \{(x^1, r^1), \dots, (x^M, r^M)\}$ ,  $r^m \in \{-1, +1\}$

$\Omega_T$  ist absolut trennbar, falls  $w := (w_1, \dots, w_I)^T$  und  $w_0$  existieren, mit

$$\begin{aligned} w^T x^m + w_0 &\geq 1, & \text{falls } r^m = 1 \\ w^T x^m + w_0 &\leq -1, & \text{falls } r^m = -1 \end{aligned}$$

# Trennband maximaler Breite

Zusammenfassung beider Fälle:

$$r^m(w^T x^m + w_0) \geq 1, \quad \forall m \in \{1, \dots, M\} \quad (1)$$

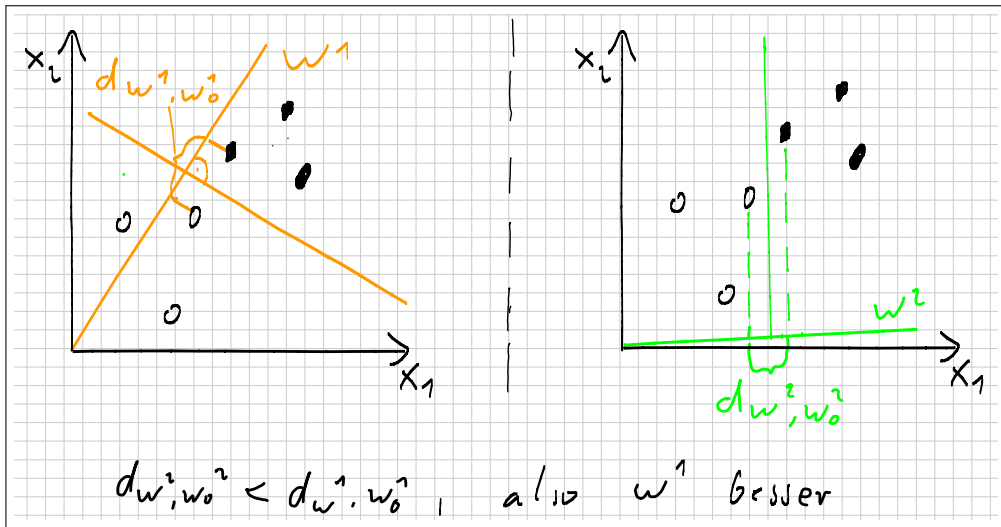
Optimale separierende Hyperebene ist definiert durch:

$$w^{*,T} x + w_0^* = 0,$$

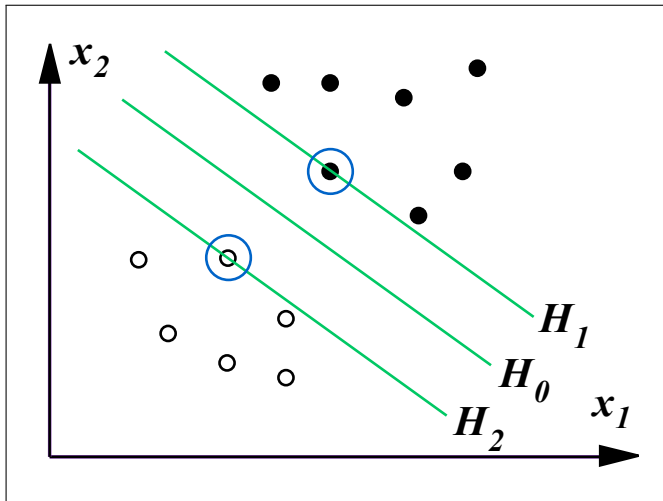
wobei  $\frac{w^*}{\|w^*\|}$  diejenige Richtung ist, bei welcher nach Projektion der Trainingsvektoren der Abstand zwischen den resultierenden Randpunkten der Klassen maximal ist. Genannter Abstand ist definiert durch:

$$d_{w,w_0} := \min_{\{(x,r=1)\}} \frac{w^T x}{\|w^T\|} - \max_{\{(x,r=-1)\}} \frac{w^T x}{\|w^T\|}$$

# Trennband maximaler Breite



# Trennband maximaler Breite



# Trennband maximaler Breite

Für einen Trainingsvektor auf  $H_1$  gilt:  $w^T x + w_0 = 1$

$\Rightarrow$  orthogonaler Abstand zum Ursprung:  $\frac{(1-w_0)}{\|w\|}$

Für einen Trainingsvektor auf  $H_2$  gilt:  $w^T x + w_0 = -1$

$\Rightarrow$  orthogonaler Abstand zum Ursprung:  $\frac{(-1-w_0)}{\|w\|}$

Also  $d_{w,w_0} := \frac{2}{\|w\|}$ .

Gesucht: Optimale separierende Hyperebene durch Minimierung von  $w^T w$  unter der Nebenbedingung gemäß Gleichung (1).



# Quadratische Optimierungsaufgabe

Standardansatz mit der Lagrange-Funktion:

$$f_{Lagr}(w, w_0, B) :=$$

$$\frac{1}{2}w^T w - \sum_{m=1}^M \beta^m (r^m (w^T x^m + w_0) - 1) \quad (2)$$

Sattelpunktfunktion  $f_{Lagr}$  in  $M + I + 1$  Parametern.

Lagrange-Multiplikatoren  $B := (\beta^1, \dots, \beta^M)^T$ , nicht-negativ.

# Quadratische Optimierungsaufgabe

Ausmultiplizieren von Gleichung (2) liefert:

$$\begin{aligned} f_{Lagr}(w, w_0, B) &:= \frac{1}{2} w^T w + \sum_{m=1}^M \beta^m - \sum_{m=1}^M \beta^m r^m w_0 - \\ &\quad \sum_{m=1}^M \beta^m r^m w^T x^m \end{aligned} \quad (3)$$

# Quadratische Optimierungsaufgabe

Satz aus der Theorie der Quadratischen Optimierung:

Sattelpunkt ist Lösung der Optimierungsaufgabe.

Minimierung bezüglich  $w, w_0$  und Maximierung bezüglich  $B$ .

Lösung erhält man prinzipiell durch Anwendung einer mathematischen Standardmethode (z.B. in MATLAB).

# Hinweis auf die optimale separierende Hyperebene

Bestimmung von  $w^*$ :

$$\begin{aligned}\nabla_w f_{Lagr}(w, w_0, B) &= w - \sum_{m=1}^M \beta^m r^m x^m \stackrel{!}{=} \vec{0} \\ \Rightarrow \quad w &= \sum_{m=1}^M \beta^m r^m x^m\end{aligned}\tag{4}$$

Hinweis: Optimale separierende Hyperebene basiert auf Linearkombination von denjenigen Trainingsvektoren mit  $\beta^m > 0$ .

Bestimmung von  $w_0^*$ :

Siehe später in diesem Unterkapitel.

# Äquivalente quadratische Optimierungsaufgabe

Nachfolgend wird eine äquivalente quadratische Optimierungsaufgabe hergeleitet. Dabei zeigt sich, daß

- nur diejenigen Trainingsvektoren von  $\Omega_T$  eine Rolle spielen, die auf  $\mathbf{H}_1$  und  $\mathbf{H}_2$  liegen, welche somit als Support-Vektoren bezeichnet werden;
- beim Training und bei der Anwendung die Input-Vektoren aus  $\Omega$  immer nur paarweise als Skalarprodukte vorkommen, auf dessen Grundlage auch nicht-lineare Diskriminanzfunktionen lernbar sind (siehe später in Unterkapitel 8.2).

# Äquivalente quadratische Optimierungsaufgabe

In der Lösung der quadratischen Optimierung gilt insbesondere:

$$\frac{\partial f_{Lagr}(w, w_0, B)}{\partial w_0} = \sum_{m=1}^M \beta^m r^m \stackrel{!}{=} 0 \quad (5)$$

Einsetzen Gleichungen (5) und (4) in Gleichung (3) liefert:

$$\begin{aligned} f_{Lagr}(w, w_0, B) &= \sum_{m=1}^M \beta^m - \frac{1}{2} \left( \underbrace{\sum_{m=1}^M \beta^m r^m x^m}_w \right)^T \left( \underbrace{\sum_{l=1}^M \beta^l r^l x^l}_w \right) \\ &= B^T U - \frac{1}{2} B^T D B =: f_{\text{Äqui}}(B) \end{aligned} \quad (6)$$

mit  $U := (1, \dots, 1)^T$ ,  $D := (D_{ml})_{m=1 \dots M, l=1 \dots M}$ ,  $D_{ml} := r^m r^l x^{m,T} x^l$

# Äquivalente quadratische Optimierungsaufgabe

Äquivalente Optimierungsaufgabe: Maximierung von  $f_{\text{Äqui}}(\mathbf{B})$  bezüglich  $\mathbf{B}$  unter den Nebenbedingungen  $\beta^m \geq 0, \forall m \in \{1, \dots, M\}$ .

Am Sattelpunkt  $\mathbf{w}^*, \mathbf{w}_0^*, \mathbf{B}^*$  gilt (nach Kuhn/Tucker):

$$\beta^{*,m} \underbrace{(\mathbf{r}^m(\mathbf{w}^{*,T} \mathbf{x}^m + \mathbf{w}_0^*) - 1)}_{(7')} = 0, \quad \forall m \in \{1, \dots, M\} \quad (7)$$

Falls  $\beta^{*,m} \neq 0$ , dann muß Ausdruck (7') gleich null sein, d.h.  $\mathbf{x}^m$  liegt auf  $H_1$  oder  $H_2$ .

Falls  $\beta^{*,m} = 0$ , dann kann  $\mathbf{x}^m$  dahinter liegen.

# Berechnung der optimalen separierenden Hyperebene

Berechnung von  $\mathbf{w}^*$  basierend auf Gleichung (4):

$$\mathbf{w}^* := \sum_{m=1}^M \beta^{*,m} r^m \mathbf{x}^m$$

Hinweis: zur Repräsentation von  $\mathbf{w}^*$  tragen nur die auf  $H_1$  oder  $H_2$  liegenden Trainingsvektoren (Support-Vektoren) bei.

Berechnung von  $\mathbf{w}_0^*$ :

Sei  $\mathbf{x}^m$  ein Support-Vektor, z.B. derjenigen Klasse mit  $r^m = 1$ .

Dann gilt  $\mathbf{w}^{*,T} \mathbf{x}^m + w_0 = 1 \Rightarrow w_0^* := 1 - \mathbf{w}^{*,T} \mathbf{x}^m$



## 8.2 SVM mit nicht-linearer Diskriminanzfunktion

### Überblick:

- Skalarprodukte von Vektoren
- Transformation des Eingaberaumes
- Implizite Raumtransformation via Kernfunktion
- Beispiele von Kernfunktionen

# Skalarprodukte von Vektoren

- In zentraler Optimierungsfunktion (6) tauchen die Trainingsvektoren nur paarweise als Skalarprodukte auf.
- Auch in der Anwendung der Diskriminanzfunktion, d.h. bei Klassifikation eines Eingabevektors  $x$ , taucht dieser Vektor nur in Skalarprodukten mit Trainingsvektoren  $x^m$  auf.

Diskriminanzfunktion:  $f_{Disk}(x) := \text{sign}(w^{*,T}x + w_0^*)$

$$\text{mit } w^{*,T}x := \sum_{m=1}^M \beta^{*,m} r^m x^{m,T}x$$

- Die Skalarprodukte werden im Folgenden durch sogenannte Kernfunktionen ersetzt. Dadurch sind auch nicht-lineare Diskriminanzfunktionen erlernbar.

# Transformation des Eingaberaumes

Transformation eines  $I_1$ -dimensionalen Eingabevektors in einen  $I_2$ -dimensionalen Vektor eines anderen Raumes.

$$f_{Tran} : \mathbb{R}^{I_1} \rightarrow \mathbb{R}^{I_2}$$

$$f_{Tran}(x) := \begin{pmatrix} f_{Tran}^1(x) \\ \vdots \\ f_{Tran}^{I_2}(x) \end{pmatrix}$$

# Transformation des Eingaberaumes

Konstruktion des Gewichtsvektors im transformierten Raum:

$$\mathbf{w}^* := \sum_{m=1}^M \beta^{*,m} \mathbf{r}^m \mathbf{f}_{Tran}(\mathbf{x}^m)$$

Nur diejenigen Summanden mit  $\beta^{*,m} > 0$  sind relevant.

Diskriminanzfunktion im transformierten Raum:

$$f_{Disk}(\mathbf{x}) := \text{sign}(\mathbf{w}^{*,T} \mathbf{f}_{Tran}(\mathbf{x}) + w_0^*)$$

# Transformation des Eingaberaumes

Paarweises Auftreten von  $f_{Tran}$  bei der Anwendung:

$$\mathbf{w}^{*,T} \mathbf{f}_{Tran}(\mathbf{x}) := \sum_{m=1}^M \beta^{*,m} \mathbf{r}^m \mathbf{f}_{Tran}(\mathbf{x}^m)^T \mathbf{f}_{Tran}(\mathbf{x})$$

Paarweises Auftreten von  $f_{Tran}$  bei Optimierung von  $B$ :

$$D_{m\ell} := \mathbf{r}^m \mathbf{r}^\ell \mathbf{f}_{Tran}(\mathbf{x}^m)^T \mathbf{f}_{Tran}(\mathbf{x}^\ell)$$

Es gibt Kombinationen aus bestimmter Transformationsfunktion  $f_{Tran}$  und einer zugehörigen, sogenannten Kernfunktion  $f_{Kern}$ , für die gilt:

$$\mathbf{f}_{Tran}(\mathbf{x}^i)^T \mathbf{f}_{Tran}(\mathbf{x}^j) \stackrel{!}{=} f_{Kern}(\mathbf{x}^i, \mathbf{x}^j), \quad \forall \mathbf{x}^i, \mathbf{x}^j$$

# Implizite Raumtransformation via Kernfunktion

Beispiel:

$$f_{Tran} : \mathbb{R}^2 \rightarrow \mathbb{R}^3, f_{Kern} : \mathbb{R}^2 \times \mathbb{R}^2 \rightarrow \mathbb{R}$$

$$f_{Tran}(x) := \begin{pmatrix} (x_1)^2 \\ \sqrt{2}x_1x_2 \\ (x_2)^2 \end{pmatrix}$$

$$f_{Kern}(x^i, x^j) := (x^{i,T}x^j)^2$$

$$\Rightarrow f_{Tran}(x^i)^T f_{Tran}(x^j) \stackrel{!}{=} f_{Kern}(x^i, x^j), \quad \forall x^i, x^j$$

# Implizite Raumtransformation via Kernfunktion

Beim Lernen muß dann nur Kernfunktion  $f_{Kern}$  herangezogen werden:

$$D_{m\ell} := r^m r^\ell f_{Kern}(x^m, x^\ell)$$

Diskriminanzfunktion:

$$f_{Disk}(x) := \text{sign} \left( \sum_{m=1}^M \beta^{*,m} r^m f_{Kern}(x, x^m) + w_0^* \right)$$

Nur diejenigen Summanden mit  $\beta^{*,m} > 0$  sind relevant.

In Artikeln von *Mercer/Courant* findet man Zusammenhänge von Transformations- und Kernfunktionen.

# Beispiele von Kernfunktionen

Gaußfunktion:

$$f_{Kern}(x^i, x^j) := e^{-\frac{\|x^i - x^j\|^2}{2\sigma^2}}$$

- Lokalisierende Gauß-(Basis)funktion als Kernfunktion.
- Diskriminanzfunktion ist durch dieselbe Art von Basisfunktionen wie beim RBF-Netz definiert.
- Im Gegensatz zur Clusterung und Verwendung von Cluster-Zentren werden die Zentren der Gauß-Funktionen hier durch Support-Vektoren definiert.

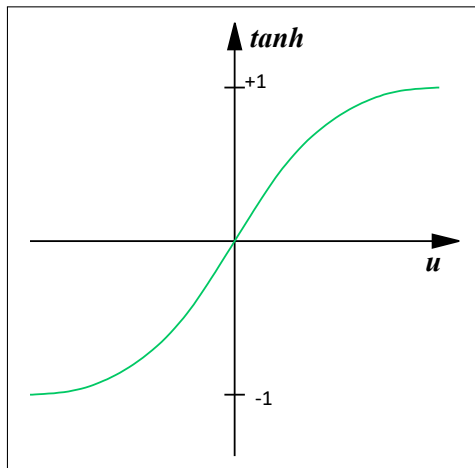


# Beispiele von Kernfunktionen

Hyperbolischer Tangens:

$$f_{Kern}(x^i, x^j) := \tanh(x^{i,T} x^j)$$

$$\tanh(u) := \frac{\sinh(u)}{\cosh(u)} = \frac{\frac{e^u - e^{-u}}{2}}{\frac{e^u + e^{-u}}{2}}$$



# Beispiele von Kernfunktionen

## Hyperbolischer Tangens (forts.):

- Nicht-Lokalisierende Sigmoid-ähnliche Kernfunktion, als kleine Gemeinsamkeit mit MLPs.
- Jedoch wird bei MLP der Sigmoid in mehreren, aufeinander folgenden Schichten jeweils auf die linearen Assoziationen von vorherigem Schicht-Output und Gewichtsvektoren für nachfolgende Knoten angewendet.
- Dagegen wird der Hyperbolischer Tangens bei SVM auf nur einer Schicht auf die linearen Assoziationen von Eingabevektor und ausgewählten Trainingsvektoren angewendet.

# Beispiele von Kernfunktionen

Polynom:  $f_{Kern}(x^i, x^j) := (x^{i,T} x^j + 1)^d$

Entspricht polynomiale Klassifikation, hier z.B mit Polynomgrad  $d$ .

- B. Schölkopf, A. Smola: Learning with Kernels - Support Vector Machines, Regularization, Optimization, and Beyond (Adaptive Computation and Machine Learning), MIT Press, Cambridge, MA, 2002, ISBN 0-262-19475-9.